

(19) World Intellectual Property Organization  
International Bureau



(43) International Publication Date  
1 November 2001 (01.11.2001)

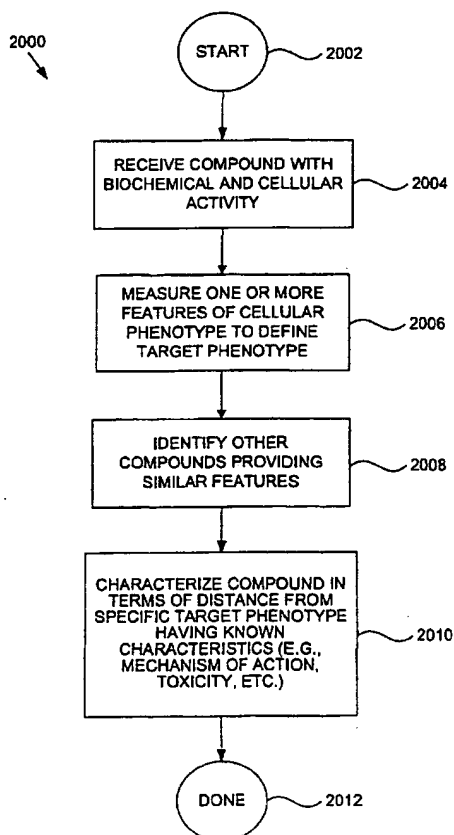
PCT

(10) International Publication Number  
**WO 01/81895 A2**

- (51) International Patent Classification<sup>7</sup>: G01N 15/00 (72) Inventors; and  
(75) Inventors/Applicants (for US only): OESTREICHER, Donald, R. [US/US]; 904 Old Town Court, Cupertino, CA 95014-4024 (US). SABRY, James, H. [CA/US]; 4305 20th Street, San Francisco, CA 94114 (US). ADAMS, Cynthia, L. [US/US]; 2409 Cedar Street, Berkeley, CA 94708 (US). VAISBERG, Eugen, A. [US/US]; 647 Pegasus Lane, Foster City, CA 94404 (US). CROMPTON, Anne, M. [US/US]; 2 Bellaire Place, San Francisco, CA 94133 (US).
- (21) International Application Number: PCT/US01/13248
- (22) International Filing Date: 24 April 2001 (24.04.2001)
- (25) Filing Language: English
- (26) Publication Language: English
- (30) Priority Data:  
60/199,778 26 April 2000 (26.04.2000) US  
09/790,214 20 February 2001 (20.02.2001) US
- (74) Agent: WEAVER, Jeffrey, K.; Beyer Weaver & Thomas, LLP, P.O. Box 778, Berkeley, CA 94704-0778 (US).
- (81) Designated States (national): AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, BZ, CA, CH, CN, CR, CU, CZ, DE, DK, DM, DZ, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ,
- (71) Applicant (for all designated States except US): CYTOKINETICS, INC. [US/US]; 280 East Grand Avenue, Suite 2, South San Francisco, CA 94080 (US).

[Continued on next page]

(54) Title: METHOD AND APPARATUS FOR PREDICTIVE CELLULAR BIOINFORMATICS



(57) Abstract: Techniques for using information technology in therapeutics or drug discovery. In an exemplary embodiment, techniques for determining information about the properties of substances based upon information about structure of living or non-living cells exposed to substances are provided. A method according to the present invention enables researchers and/or scientists to identify promising candidates in the search for new and better medicines or treatments using, for example, a multiple biological descriptors derived from a single cell component or marker. The method employs image analysis to extract a plurality of features (e.g., cell size, distance between cells, cell population, cell type) from an image acquisition device into the database.

WO 01/81895 A2

BEST AVAILABLE COPY



NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM,  
TR, TT, TZ, UA, UG, US, UZ, VN, YU, ZA, ZW.

(84) **Designated States (regional):** ARIPO patent (GH, GM, KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE, TR), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GW, ML, MR, NE, SN, TD, TG).

**Published:**

— without international search report and to be republished upon receipt of that report

*For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.*

## **METHOD AND APPARATUS FOR PREDICTIVE CELLULAR BIOINFORMATICS**

### **COPYRIGHT NOTICE**

A portion of the disclosure of this patent document contains material which is subject to copyright protection. The copyright owner has no objection to the facsimile reproduction by anyone of the patent document or the patent disclosure as it appears in the Patent and Trademark Office patent file or records, but otherwise reserves all copyright rights whatsoever. The present description includes examples of computer codes, which may be used to implement aspects of the present invention. Assignee of the present invention reserves all rights with respect to these codes and provides notice herein. Notice is hereby given © Cytokinetics, Inc. 1999, 2000, 2001.

### **BACKGROUND OF THE INVENTION**

The present invention provides techniques for information management using a database platform. More particularly, the present invention provides a system including computer code that couples to a database device. The system provides for image capturing of living, dead, or fixed cells or cell fractions used to identify information about substances used on the cells or information about the cells themselves. Accordingly, the present invention can enable researchers and scientists to identify promising candidates in the search for new and better medicines, for example, in drug discovery and development. The principles enumerated herein may, with equal facility, be applied to other applications, including but not limited to use in environmental applications such as determining chemical toxicities and other non-pharmaceutical toxicology uses.

For a long time, researchers in the pharmaceutical field have sought for better ways of searching for substances possessing properties that make them suitable as medicines. In the early days, researchers generally relied upon extracts from plants, dyes, and microbiological extracts for such substances. Examples of such substances include the pain reliever aspirin, the anti-cancer drug paclitaxel (brand name Taxol<sup>TM</sup>), and the heart medication called digoxin. The number of useful medicines has generally been limited.

Purified substances having desirable bio-active properties are also often difficult to discover. Advances in traditional organic chemistry and more recently the rapid chemical synthesis methods often referred to as combinatorial chemistry have increased the number of compounds that researchers test for biological activity. Originally, substances were often initially tested on animals or humans to determine their biological activity. While results from such tests may identify a good drug candidate, they are often time consuming and costly, thus a limited number of substances can be tested. Therefore, pharmaceutical companies have turned to testing their ever-increasing libraries of substances against isolated proteins (drug targets) in biochemical assays that can be carried out at high throughput and low cost. It should be noted that the substances need to be tested in numerous protein tests, each customized for a particular drug target. Therefore, although each protein test may be run at a high-throughput, the design of multiple protein tests can be time-consuming. Substances deemed promising based on results from the protein tests are then tested in lower throughput cellular and animal tests.

There have been some attempts to use image acquisition techniques to screen a large number of substances based upon biological cell information. One such attempt is described in International Application No. WO 98/38490 in the names of Dunlay, et al. Dunlay et al. generally describes a conventional image acquisition system. This conventional system collects and saves images based on certain criteria that are predefined, not on a fixed area of an imaging surface. Additionally, the conventional system has poor lighting design, which makes image processing for multiple cells difficult. Furthermore, the conventional system is not designed for capturing, populating and utilizing a large database design. The conventional system is designed for customized cellular assays, not as a tool for generation of a cellular informatics database. Without such database capabilities the conventional system cannot be used for screening, analyzing, and comparing large quantities of cells from multiple experiments on multiple days in a predictive, efficient and cost effective manner.

What is needed is a rapid assay to assess the activity of compounds against multiple drug targets simultaneously in a cellular context. What is also needed are techniques for finding the effects of substances on cell function based upon searching and analyzing cellular information.

### SUMMARY OF THE INVENTION

According to at least one embodiment of the present invention, techniques for determining information about effects of potential substances on cells are provided. In another exemplary embodiment, the present invention provides a novel system including hardware, computer codes, user interfaces, and a database for acquiring, storing and retrieving cellular and substance information. The cells can include living, dead, or fixed cells or fractions of cells. The present invention enables, *inter alia*, researchers and/or scientists to identify promising candidates in the search for new and better medicines or treatments using, for example, a cellular informatics database.

According to the present invention, a computer program for identification and verification of biological properties of substances can include code that causes a sample of a substance to be administered to a cell. The code determines one or more features for two or more cell components, or markers, in the presence of the substance. The code can form one or more descriptors from the features. Descriptors can be formed by combining features of two or more cell components as identified using the markers. The code can then search one or more descriptors obtained from prior administered substances upon cells in order to locate descriptors having a relationship to the descriptors noted for the substance under study. The code predicts properties of the administered substance based upon the properties of the prior administered substances using the relationship between the descriptors. The code can provide for identifying properties of substances based upon effects on cell characteristics. Candidate drug mechanisms of action, potency, specificity, pharmacodynamic, and pharmacokinetic parameters, toxicity, and the like can be used as substance properties.

In a specific embodiment, the present invention provides a system for acquiring knowledge from cellular information. The system has a database comprising a database management module ("DBMS"). The system also has a variety of other modules, including a population module that is coupled to the DBMS and serves to categorize and store a plurality of features (including but not limited to cell size, distance between cells, cell population, as well as sub-cellular features such as organelle location, protein location and sub-cellular constituent location and movement) from an image acquisition device into the database. The system has a translation module coupled to the DBMS for defining a descriptor from a set of selected features from the plurality of features. In a specific embodiment, the descriptor is for a known or unknown compound,

e.g., drug. A prediction module is coupled to the DBMS for selecting one of a plurality of a descriptors from known and unknown compounds from the database based upon a selected descriptor from a selected compound. The selected compound may be one that is useful for treatment of human beings or the like.

In a specific embodiment, the present invention provides a system for populating a database with cellular information. The system includes a cell holder (e.g., multi-well plate, chip, microfluidic assembly, or other cell chamber) comprising a plurality of sites in a spatial orientation. Each of the sites is capable of holding a plurality of cells to be imaged. Note – the light guide is one embodiment, but we don't want to be limited to it.

According to one embodiment, the present system also has an illumination apparatus including a liquid light guide operably coupled to the imaging device for highlighting the plurality of cells in a relatively even spatial manner for image capturing and measurement purposes. Still further, the liquid light guide allows sub-elements (e.g., filter, lamp) of the illumination apparatus to be placed at a remote location to prevent mechanical interference of the cell holder during image capturing. Alternative lighting methodologies may, with equal facility, be implemented.

The system also has an image-capturing device (e.g., charge coupled device camera, translation stage, shutter, microscope, software, shutter control) coupled to a computing device (e.g., computer, network computer, work station, analog computing device, on-board image-processor, and laptop). The image-capturing device is adapted to capture at least one image in at least one of the plurality of sites. In some embodiments, multiple images can be captured, where each image represents a different cell component (or portion). The image-capturing device can be adapted to convert the image into a digital representation, which highlights the feature or features of the one site.

A database storage device (e.g., relational database, object oriented database, mixed object oriented database) includes a database management element. The database is coupled to the image capturing device. In a specific embodiment, the present system includes modules for feature extraction, generation of descriptions, and data preparation and analysis.

In a specific embodiment, the present invention provides a novel system for determining an effect of a manipulation of a cell using one or more image frames.

The system has a plate comprising a plurality of sites in a spatial orientation. Each of the sites is capable of holding a plurality of cells to be imaged. The system also has an image capturing device to capture a plurality of images of at least one site from the plurality of sites. The image capturing device is coupled to the computing device. The system also has an image processing device to combine the plurality of images of at least one site or plurality of sites. The image processing device is operably coupled to the plate. An image processing device is also included. The image processing device can be adapted to form a digitized representation of the plurality of images from the site or plurality of sites. Furthermore, the system has a database storage device comprising a database management element. The database can be adapted to retrieve the descriptor or descriptors of the plurality of features from the computing processing device and storing them in a selected manner.

In a specific embodiment, the present invention provides a system for capturing cellular information. The system also has an image acquisition system comprising a charged coupled device camera adapted to capture an image of a plurality of manipulated cells in various stages of the cell cycle. The stages of the cell cycle are currently understood to include interphase, G0 phase, G1 phase, S phase, G2 phase, M phase, prophase, prometaphase, metaphase, anaphase, and telophase. The principles of the present invention specifically contemplate the application thereof on additional cell cycle stages when and if they are identified.

An optical source is coupled to the image acquisition system for highlighting the plurality of manipulated cells in the various stages of the cell cycle. The illumination apparatus provides for an acquisition of the image of the plurality of manipulated cells. In a specific embodiment, the illumination apparatus has a liquid light guide coupled to a light source at a remote location.

A variety of user interfaces are utile for accessing the several features of the present invention. Those having ordinary skill in the art will appreciate that different user interfaces may be required to support different research scenarios. The present invention specifically contemplates the utilization of a wide variety of user interfaces.

Numerous benefits are achieved by way of the present invention over conventional techniques. The present invention can provide techniques for predictive cellular bioinformatics that can streamline a number of important decisions made in the drug discovery industry. The present invention can be implemented using off the shelf

hardware including databases. In other aspects, the present invention can find useful information about substances as well as cells or portions of cells. Furthermore, the present invention can acquire more than one feature using more than one manipulation. Moreover, the present invention can provide information about a wide variety of cellular information that is not conventionally available. This information includes information about different cell components, e.g., nuclei and Golgi apparatus. Still further, the present invention provides an automated or semi-automated technique for acquiring images and populating a database. The present database can be combined with others such as genomics, and the like. Moreover, the present invention can be implemented to predict, *inter alia*, a mechanism of action, toxicity, target validation, and pre-clinical disease model.

A further understanding of the nature and advantages of the invention herein may be realized by reference to the remaining sections of the specification and the attached drawings.



### BRIEF DESCRIPTION OF THE DRAWING

For more complete understanding of the present invention, reference is made to the accompanying Drawing in the following Detailed Description of the Invention. In the drawing:

Fig. 1 is a simplified system diagram according to an embodiment according to the present invention;

Figs. 1A-1B are more detailed diagrams of database systems according to embodiments of the present invention;

Fig. 2 is a simplified block diagram according to an alternative embodiment according to the present invention;

Figs. 3-6 are simplified diagrams of system elements according to embodiments of the present invention;

Figs. 7A-7K illustrate representative block diagrams of simplified process steps in a particular embodiment according to the present invention;

Fig. 8A-8F illustrate representative quantified descriptors of effects of manipulations on images of cells in a particular experiment;

Fig. 9 illustrates example images for different types of morphologies in a particular experiment;

Fig. 10 illustrates a distribution of various morphologies in a cell population responsive to drug concentration in a particular experiment;

Fig. 11 illustrates a graph of quantified features of effects of manipulations on cells in a particular experiment;

Fig. 12 illustrates effects of external agents on cells in a particular experiment;

Fig. 13 illustrates 4 panels for each marker for a plurality of A549 cells in a particular experiment;

Fig. 14 illustrates 4 panels for each marker for a plurality of OVCAR-3 cells in a particular experiment;

Fig. 15 illustrates 4 panels for each marker for a plurality of OVCAR-3 cells at 20x in a particular experiment;

Fig. 16 illustrates 4 panels for each marker for a plurality of OVCAR-3 cells at 40x in a particular experiment;

Fig. 17 illustrates a representative input for a morphometric analysis program in a particular embodiment according to the present invention; and

Figs. 18-19 illustrate examples of the generation of pseudo-sequences and clustering in a particular embodiment according to the present invention.

Fig. 20 is a block diagram for a first research scenario;

Fig. 21 is a block diagram for a second research scenario; and

Fig. 22 is a block diagram for a third research scenario.

Fig. 23 is a process flow diagram that illustrates a process employing three separate biologically relevant measurements made from a single component/marker shown in a cell image.

Reference numbers refer to the same or equivalent parts of the invention throughout the several figures of the Drawing.

### DETAILED DESCRIPTION OF THE INVENTION

According to the present invention, techniques for determining information about manipulated cells or substances based upon living, fixed, or dead cell structures or portions of cells are provided. In an exemplary embodiment, the present invention provides a novel system including computer codes coupled to a database and user interfaces for acquiring, storing and retrieving such information. Other embodiments provide a novel image capturing system for providing digitized representations of live and dead cell structures or the like.

Fig. 1 is a simplified system diagram 10 of a cellular knowledge-based system according to an embodiment of the present invention. This diagram is merely an example and should not limit the scope of the claims herein. One of ordinary skill in the art would recognize other variations, modifications, and alternatives. The present system 10 includes a variety of elements such as a computing device 13, which is coupled to an image processor 15 and is coupled to a database 21. The image processor receives information from an image capturing device 17, which image processor and image capturing device are collectively referred to as the imaging system herein. The image capturing device obtains information from a plate 19, which includes a plurality of sites for cells. These cells can be biological cells that are living, fixed, dead, cell fractions, cells in a tissue, and the like. The computing device retrieves the information, which has been digitized, from the image processing device and stores such information into the database. A user interface device 11, which can be a personal computer, a work station, a network computer, a personal digital assistant, or the like, is coupled to the computing device.

Fig. 1A is a simplified diagram of a database system 1000 according to an embodiment of the present invention. This diagram is merely an example and should not limit the scope of the claims herein. One of ordinary skill in the art would recognize many other variations, modifications, and alternatives. Database system 1000 includes a variety of techniques for processing images from biological cells, e.g., fixed, living, and dead cells, and cell portions. As shown, images are acquired 1001. These images can be from a single frame or multiple frames. As merely an example, an image processing system may analyze such images. One example of such an image processing system is described below, but should not be construed as limiting certain claims.

In a specific embodiment, cell samples are manipulated using a compound (e.g., substance, drug). The cell samples are imaged for a simple portion or portions, e.g., manipulated cell substructure, manipulated spatial feature of cell, cell density. Image processing techniques are used to extract 1003 the feature or features from the image or images. The features can be an independent or a dependent set of cell characteristics (which may be predominately visual) including, for example, count, area, perimeter, length, breadth, fiber length, fiber breadth, shape factor, elliptical form factor, inner radius, outer radius, mean radius, equivalent radius, equivalent sphere volume, equivalent prolate volume, equivalent oblate volume, equivalent sphere surface, average intensity, total intensity, optical density, radial dispersion, texture difference, and others. Each of these features corresponds to a similar manipulation by a compound. Each manipulation forms a new set of features, which are identifiable to the compound. Once each set of features has been extracted, the feature set is populated 1004 into a database 1012. Accordingly, the database includes many sets of features, where each set corresponds to a different manipulation for a selected cell. Each set of features corresponding to a manipulation provides a descriptor 1009, which is also stored 1019 in the database. The descriptor is a "finger print" including each feature for the manipulation. Each descriptor may be unique, or may have similarities to other descriptors or may even be the same as other descriptors for known and unknown manipulations.

The present system retrieves features, which we define as simple features herein, and forms composite features 1007 from them. More than one feature can be combined in a variety of different ways to form these composite features. In particular, the composite feature can be any function or combination of a simple feature and other composite features. The function can be algebraic, logical, sinusoidal, logarithmic, linear, hyperbolic, statistical, and the like. Alternatively, more than one simple feature can be combined in a functional manner (e.g., arithmetic, algebraic). As merely an example, the composite feature equals a sum of feature 1 and feature 2, where these features correspond to the same manipulation. Alternatively, the composite feature equals feature 1 divided by feature 2. Alternatively, the composite feature equals feature 1 minus feature 2. Alternatively, the composite feature equals a constant times feature 1 plus feature 2. Of course, there are many ways that the composite feature can be defined. The present system also stores 1017 these features in the database. The

composite features can also be further combined with simple features. Once these features are defined as descriptors, they are stored 1019 in the database.

Fig. 1B is a simplified diagram of a database system engine 2000 according to an embodiment of the present invention. This diagram is merely an example and should not limit the scope of the claims herein. One of ordinary skill in the art would recognize many other variations, modifications, and alternatives. The engine can be implemented into the present database for populating, searching, and predicting compound or cell characteristics. As merely an example, engine 2001 includes an input/output module 2008. The input/output module is used to input and output information from the database. The information includes, among others, a plurality of feature sets, which correspond to many manipulations. Additionally, the information includes descriptors, which each corresponds to a set of features from the manipulation. The database also has a population module, which is used to configure the features based upon an entity relationship, which has been predetermined.

The database engine also has other modules. In particular, the database has a transcription module, which transfers a preselected set of features and creates a descriptor from them. The transcription module can be used to take a known compound, which has features, to transcribe them into a descriptor. Alternatively, the transcription module can be used to take an unknown compound, which has features, to transcribe them into a descriptor. These descriptors are provided into the database for subsequent use. Finally, the database engine has a prediction module, which can be used to potentially predict a property (e.g., mechanism of action) of an unknown compound. Here, the unknown compound is provided with a descriptor, but the property of the compound is unknown. In one embodiment, the prediction module compares a descriptor of an unknown compound with the many descriptors of known compounds, which were in the populated database. Depending upon the matching criteria, the prediction module will attempt to uncover one or more descriptors of known compounds. Once the prediction module finds the descriptors of the known compounds based upon the descriptor for the unknown compound, it identifies a potential property of such unknown compound for analysis and review. Here, it is believed that certain features of the known compound, which are similar to those features of the unknown compound may uncover a property to the unknown compound. Details of the present software engine are described more fully below.

Fig. 2 is a simplified block diagram 20 of a cellular knowledge-based system according to an alternative embodiment of the present invention. This diagram is merely an example and should not limit the scope of the claims herein. One of ordinary skill in the art would recognize other variations, modifications, and alternatives. Like reference numerals are used in the present diagram as the previous diagram for easy cross-referencing, but are not intended to be limiting in any manner. The present diagram 20 includes a variety of elements such as a processor 13 or computing device coupled to a database 11. The processor can be used for retrieving and storing information from the database. The system also includes a plurality of system elements, such as a cleaner 23, a dispenser 25, and an image capturing system 27, which are also coupled to the database in some embodiments. These elements can be coupled to each other through a network or the like. As merely an example, the network can be a NetWare<sup>TM</sup> network from Novell Corporation or an internet network or the Internet but can also be others and any combination thereof. The system also has an output device 31, which can be used to output information from the database, processor, or other system elements. Details of these elements are described more fully below in reference to the Figs.

Figs. 3-5 are simplified drawings of system elements according to embodiments of the present invention. These diagrams are merely examples and should not limit the scope of the claims herein. One of ordinary skill in the art would recognize other variations, modifications, and alternatives. As merely an example, Fig. 3 is a simplified diagram of a processor or computing device 13. The computing device 13 includes a bus 112 which interconnects major subsystems such as a central processor 114, a system memory 116 (e.g., random access memory), an input/output ("I/O") controller 118, an external device such as a display screen 124 via a display adapter 126, a keyboard 132 and a mouse 146 via an I/O controller 118, a SCSI host adapter (not shown), and a floppy disk drive 136 operative to receive a floppy disk 138.

The computing device has other features. Storage Interface 134 may act as a storage interface to a fixed disk drive 144 or a CD-ROM player 140 operative to receive a CD-ROM 142. Fixed disk 144 may be a part of computing device or may be separate and accessed through other interface systems. A network interface 148 may provide a direct connection to a remote server via a telephone link or to the Internet. Network interface 148 may also connect to a local area network ("LAN") or other

network interconnecting many computer systems. Many other devices or subsystems (not shown) may be connected in a similar manner. Also, it is not necessary for all of the devices shown in Fig. 3 to be present to practice the present invention, as discussed below. The devices and subsystems may be interconnected in different ways from that shown in Fig. 3. The operation of a computer system such as that shown in Fig. 3 is readily known in the art and is not discussed in detail in this application. Computer code to implement the present invention, may be operably disposed or stored in computer-readable storage media such as system memory 116, fixed disk 144, CD-ROM 140, or floppy disk 138. The computer code can be organized in terms of processes or modules, depending upon the application. That is, the computer code can include a prediction module, a translation module, or other modules to carry out the functionality described herein, as well as others.

Figs. 4 and 5 are simplified diagrams of an imaging system 200 according to an embodiment of the present invention. As shown, the imaging system 200 includes a variety of features such as housing 203, which holds a stage assembly 204. The stage assembly includes an x-stage movement element 206, which is along an x-direction, and a y-stage movement element 207, which is along a y-direction. The imaging system also includes a z-direction movement element, which is perpendicular to the x-y plane. The z-direction movement motor can be attached to the stage, or to the objective nosepiece by way of the microscope housing, or as an external motor between the objective and the microscope housing. The stage can align in any one of the directions to an accuracy of one micron and less, or one-half micron and less, or one-quarter micron and less, depending upon the embodiment.

The stage holds a plate 202 or cell holder, which houses one of a plurality of samples. The plate includes a spatial array 209 of process sites. Each of the process sites can include a plurality of cells and solutions depending upon the embodiment. Each of the sites can carry a sufficient amount of solution to prevent substantial evaporation of the sample during processing in some embodiments. In embodiments for large scale analysis, the plate includes at least 96 sites, or more than or equal to 384 sites, or more than or equal to 1,536 sites. The plate bottom is transparent and thin, which allows light to pass through the sample. Additionally, the plate is made of a suitable chemical resistant material. As merely an example, the plate can be either a 96, or 384, or 1536 or other formats from places such as Becton Dickinson of Franklin Lakes, NJ, or Corning

Science Products of Corning, NY. In a preferred embodiment, the plate is a Corning Costar black-walled 96 well plate catalog #3904 from Corning Science Products of Corning, NY, but should not be limited to these in some applications, but can be others.

Also shown is the condenser for the microscope 201, which can be used to collect phase, DIC, or bright field images of the cells. Images resulting from the illumination of the samples to fluorescence, phase, DIC, or bright field techniques are collected using an image capturing device 208, which captures an image or images of cells from the plate. In a specific embodiment, the microscope is an inverted configuration with the objectives on the bottom of the plate and the condenser disposed overlying an upper surface of the sites, while the image capturing device underlies the sites. Images captured by the imaging device, whether analogue or digital, are viewed by a monitor or other devices. The image capturing device can be any camera assembly such as a charge coupled device camera, which is known as a CCD camera, or other high resolution camera capable of capturing images from the sites. In a specific embodiment, the camera is an interline CCD camera which does not require an external shutter.

In a specific embodiment, the present imaging system can be any suitable unit that is flexible for automated image collection using multi-well plastic plates. The imaging system also should be adapted to collect high-resolution images of cells on plastic or glass plates, cell growth chambers, or coverslips. The system also can be used for imaging multiple cell markers in multiple imaging conditions. To accomplish this, the microscope system has a variety of elements such as a light source, a motorized excitation filter wheel and shutter, x-y-z-motorized stage, excitation and emission filters, Fluor phase and DIC objectives, motorized objective nosepiece, dichroic filters, motorized dichroic filter cubes, phase and DIC rings and prisms, CCD camera, and software control. As merely an example, the present imaging system can have components such as those listed in the Table below.

DESCRIPTION	MAKER	MODEL
Microscope	Zeiss	100M
(x-y) motorized stage	Prior	
Xenon lamp	Sutter	Lambda
Filter wheel	Sutter	Lambda-10



Microtitre Plate holder	Prior	500-H223R
Isolation Table	Kinetic Systems	9101-24-85
Objective Spacers	Polytec PI	P-721.90
Camera	Hamamatsu	C47-95
Computer	IBM	IntelliStation
Software	Metamorph	v.4
Objectives	Zeiss	Achroplan 10x/0.25 LD-Achroplan 20x/0.4 LD-Achroplan 40x/0.6

Table: Image Acquisition System Elements

In a specific embodiment, the present system has the following capabilities, which are not intended to be limiting.

#### Image acquisition

1) Ability to automatically acquire multi-wavelength images from multiple sites on one multi-well plate, to sequentially name image files, and to log any imaging parameter information with image files.

2) Ability to link images with a larger database/spreadsheet of information.

3) Ability to automatically collect multiple plates by interfacing the imaging system with a robotic arm.

#### X-Y control

1) Ability to place 96, 384, or 1536 well plates onto microscope stage and move to each well sequentially.

2) Ability to return to each well and collect another round of images (multi-site time-lapse) or ability to collect rapid time-lapse information at each well (time-lapse of many wells).

3) Ability to collect a low magnification image, automatically determine features which may be of interest, automatically change the objective to a higher magnification, and collect high magnification images of a fixed number of those identified cells in the sample.

4) Ability to collect multiple frames in each site.

Z control

1. Ability to auto-focus with substantially minimal damage to biological specimen or fluorophore.

2. Ability to auto-focus rapidly.

The present embodiment of the imaging system is shown by way of Figs. 5A and 5B. These diagrams are merely examples and should not limit the scope of the claims herein. One of ordinary skill in the art would recognize other variations, modifications, and alternatives. The present imaging system 40 includes a variety of elements such as a microscope 41, which is preferably an epi-fluorescent microscope, but can be confocal, multiphoton, or hybrid types. The microscope includes elements 41A, the motorized Z-axis; 41B, the motorized dichroic filter cube holder; and 41C, the motorized objective nosepiece. In one embodiment, the microscope is a Model 100M made by Zeiss. The microscope communicates to computer 51 through control lines 73, 75, and 76. The imaging system also has camera 50 coupled to controller 50A and computing device 51, which oversees and controls operations of the elements of the imaging system.

The present microscope includes drivers for spatially moving a stage in two dimensions, including an x-direction, a y-direction, and moving the objective nosepiece in a z-direction in a Cartesian coordinate system. The z-direction movement is provided using a fast z-motor, which can make z-direction adjustments within a predetermined time. The z-direction movement generally provides for focussing of the sample to the camera. The focussing occurs within the predetermined time of preferably ten seconds and less, or five seconds and less, or one second and less, depending upon the embodiment. As merely an example, the z-motor or positioner can be a model PIFOC objective nanopositioner made by a company called Physik Instrumente of Waldbronn, Germany, but also can be others. The z-motor couples to computer 51 through line 63, which may also include a controller. Depending upon the embodiment, a second z-motor 41A connected to the computer 51 by line 73 may be used to keep the z-motor 42 in the center of its travel. Alternatively, in other embodiments the stage could be provided with a z-motor allowing for movement of the stage in the z-direction.

The present stage also includes an x-y stage 43. The x-y stage moves plate 59, e.g., 96 site, 384 site, 1536 site. The x-y stage moves plate in an x-y spatial manner. The stage has an accuracy or repeatability of about 1 micron and less, or about 2 microns and less. The stage can move in a continuous manner or a stepped manner. The stage also can move up to 30 mm/sec. or faster. The stage also can move 1 mm/sec. and less, depending upon the embodiment. The stage can also step 0.1 micron and less or 1 micron and less, as well as other spatial dimensions. The stage can be one such as a Proscan Series made by Prior Scientific of Rockland, MA but can also be others. The stage is controlled via control line 61 through controller 43A, which couples to computer 51 through control line 65.

The stage includes plate holder 44. The plate holder can hold a single plate. In other embodiments, plate holder can also hold multiple plates. The plate holder can use mechanical, electrical, fluid, vacuum and other means for holding the plate or plates. The plate holder also is sufficiently stable for securing the plate. As merely an example, the plate holder is a Model 500-H223R made by Prior Scientific of Rockland, MA. In some embodiments, the plate holder may need adjustment in the z-direction to provide for a desirable focus of a sample on a plate. In these embodiments, the plate holder is supported by spacers 45 or a plurality of stage pins, which mechanically elevate the plate holder in the z-direction. These pins are generally made of a suitable material for supporting such plate holder and also are sufficiently resistant to chemicals and the like.

In some embodiments, the entire imaging system is placed on an isolation table 49. The isolation table is disposed between the microscope and support structure. The isolation table is designed to prevent excessive vibration of the plate. The isolation table is made of a suitable material such as steel and honeycomb but can be others. The table has a thickness of about 8 inches or preferably less than about 24 inches. In one embodiment, the table is Model 9101-24-85 made by Kinetic Systems of Boston, MA.

The imaging system also has a lamp or illumination assembly 62. The lamp assembly provides for a light source (See reference letter B) to a plurality of elements in the imaging system. For easy reading, the light path is defined by the dotted lines, which are not intended to be limiting. The lamp assembly has a variety of elements such as a Xenon lamp 46. The Xenon lamp provides light at about 320 to 700

nanometers (Prefocused). The Xenon lamp is 175 or 300 Watts. As merely an example, the lamp can be a Lambda Model made by Sutter Instrument Company of Novato, CA.

Referring to Fig. 5B, the lamp assembly also has a cold mirror 58, an excitation filter wheel 48, excitation filter(s) 55, and an excitation light shutter 57. As shown, light is derived from the Xenon lamp, reflects off of the cold mirror 58, traverses through the excitation filter or filters 55, and is controlled by the excitation light shutter 57. The lamp assembly has filter wheel 48, which houses one of a plurality of filters, including excitation filters. The shutter and filter wheel are controlled via control lines 67, which are coupled to a computer 51 or other type of computing device. The control lines 67 are coupled through controller 57A (for element 57) and controller 48A (for element 48) via control line 69 to computer 51.

Preferably, light traverses from the lamp assembly through a light guide 47 to illuminate features within the plate. The light guide is suitably selected to have a flexible member, which can be used to place lamp source at a remote location away from the imaging device. The flexible member substantially keeps any vibration from the lamp assembly away from the imaging device. In some embodiments, the member is at least 1 foot away from the imaging device. The light guide is a guide, which is a flexible hose-type sleeve. The sleeve is filled with a liquid such as an aqueous solution containing chloride or phosphate. A thin layer may be formed on the inside of the sleeve. The layer can be a containing tetrafluoroethylene and hexafluoropropylene, or containing tetrafluoroethylene and perfluoromethyl vinyl ether, or tetrafluoroethylene and perfluoropropyl vinyl ether. An example of such a light guide is described in International Application No. WO/98/38537 filed February 29, 1997, and assigned to NATH, Gunther. The liquid light guide has less than about 30% transmission loss of the light at a remote location such as the imaging system.

Light is derived from the lamp assembly and directs off of filter 56, which directs the light upward. Filter 56 can be a dichroic and emission filter, as well as others. The light traverses through microscope nosepiece 41C, and traverses through objective spacers 54. An objective 53 magnifies the light toward a predetermined point on the plate 59. The objective can be, for example, made by Zeiss of Jena, Germany, as well as other companies. The objective can be one of a plurality including 1X, 10X, 20X, 40X, and others, depending upon the application. Magnification can be further expanded

or contracted by intermediate optics between the objective and the camera. Selection of filter or filters is controlled by computer 51 via control line 75.

The camera 50 captures an image of cells from plate 59. The image is obtained from light scattering off of cells or portions of cells in the plate through objective 53, through objective spacers, through filters 56, which are captured at camera 50. In this preferred embodiment, the camera is a digital camera, but can be an analogue camera. The digital camera is a CCD camera, which has 1280 by 1024 pixels, or more or less. The pixels can be 6.7 microns in dimension or more or less. The camera preferably is substantially free from an external shutter to quickly capture a plurality of images of cells from the plate. The camera is controlled via control line 71 through controller 50A, which connects to computer 51 through control line 70. The present invention can also include other types of image acquisition devices selected from at least an epifluorescence, a confocal, a total-internal reflection, a phase, a Hoffman, a bright field, a dark field, a differential interference contrast, an interference reflection, or multi-photon illumination device.

The present imaging system stores images on a high density memory device 60. The high density memory device is preferably optical, but can also be magnetic. The high density memory device can be any suitable unit that is capable of storing a plurality of images from a plurality of sites in the plate. The memory device can be a compact disk, which would generally use a compact disk burner or the like. Depending upon the embodiment, the high density memory device is used to archive the images that are captured from the camera in the imaging system. Further details of the imaging system can be found throughout the present specification, and more particularly below.

As merely an example, the present invention can be implemented using the following sequence of steps, which have been described in a journal entry form. Here, images are opened and objects are identified based on a background value that has been edited in starting image acquisition. Information is maintained in a spreadsheet or other database format, which has the following information for each object:

Image Name	Image Plane	Image Date and Time
Elapsed Time	Object #	Total area

Pixel area	Area	Hole area
Relative hole area	Standard area count	Perimeter
Length	Breadth	Fiber length
Fiber breadth	Shape factor	Ell. form factor
Inner radius	Outer radius	Mean radius
Average gray value	Total gray value	Optical density
Radial dispersion	Texture Difference Moment	EFA Harmonic 2, Semi-Major Axis
EFA Harmonic 2, Semi-Minor Axis	EFA Harmonic 2, Semi-Major Axis Angle	EFA Harmonic 2, Ellipse Area
EFA Harmonic 2, Axial Ratio	EFA Harmonic 3, Semi-Minor Axis	

After computations are done, the log file is saved. In particular, the file is saved in an appropriate place with an appropriate name.

In a specific embodiment, the present invention provides the following detailed example of journal entries, which should not limit the scope of the invention.

Set Up Sequential File Names	Interactive: user sets up prefix name and image storage directory
Open Data Log	Opens a DDE (Excel) File
Annotate Log File	Interactive: experimental information that will go into the first line of the log file of stage positions
Stage (Go to Origin)	Origin is set as the center of well A1
Stage (Move to Absolute Position)	Offset to upper left hand corner of well (1410, 1621)
Stage (Log Position)	

Stage (Scan Wells)	User picks wells to scan: runs 3x3 image collection.jnl.
--------------------	--

3X3 IMAGE COLLECTION.jnl

Stage (Scan)	Takes 9 images of well, -1600 motor steps apart from left to right 3 columns and 3 rows, runs FOCUS, COLLECT IMAGE, SAVE SEQUENTIAL FILE NAME.JNL.
--------------	---

FOCUS, COLLECT IMAGE, SAVE SEQUENTIAL FILE NAME.jnl.

Stage (Log Position)	Logs stage position of each image
ADC – Focus	Opens up the manual focusing window with whatever focus time is current set
Show Message and Wait	Interactive: user hits enter to continue when done focusing
ADC-Acquire from Digital Camera	Takes Hoechst image
Save Using Sequential File Names	
Close	Closes image window

START IMAGE ANALYSIS.jnl

Low Pass	3x3 convolution of already opened image
Low Pass	3x3

Show Region Statistics	Interactive: Show entire image statistics. Calculate background subtraction value for step 4. by: INTENSITY Average + INTENSITY Std. Dev.
Arithmetic	Interactive: User inputs subtraction value from 3. into the constant Value field
Threshold image	Creates threshold 1 unit above 0 to 4096
Integrated Morphometry – Load State	Loads Start Image Analysis.ima Classifier 100 < area < 200000
Integrated Morphometry – Measure	Interactive: Shows area summary information about all objects. The average number is used as the Standard Area in 8.
Object Standards - Set Object Standards	Interactive: User inputs average area value from 7. into Standard Area box to be used by automated IMA for all images

IMA OBJECTS.jnl

Low Pass	3x3 convolution
Low Pass	3x3 convolution
Arithmetic	This background subtraction value needs to be manually entered into this journal from the value determined in START IMAGE ANALYSIS.jnl step 3
Threshold Image	1 unit above 0



Integrated Morphometry – Load State	Hoechst.IMA Classifier 200 < area < 200000
Integrated Morphometry – Measure	Measures statistical info for all objects
Run Journal	Runs log obj and sum data.jnl

Log obj and sum data.jnl

Integrated Morphometry – Log Data	Logs object data into Sheet 1
Integrated Morphometry – Log Data	Log summary data into Sheet 2

## COLLECT AUTOMATED IMA DATA IN ONE SPREADSHEET.jnl

Run Journal	Runs OPEN OBJECT LOG DDE FILE.JNL
Loop for all Images in a Directory	Loops IMA OBJECTS.jnl
Close Summary Log	
Close Object Log	User must manually save Excel spreadsheet

OPEN OBJECT LOG DDE FILE.jnl

Open Object Log	Opens a DDE object log into sheet 1 of an Excel spreadsheet
Open Summary Log	Opens a summary log into sheet 2

## COLLECT AUTOMATED IMA DATA IN ONE SPREADSHEET 16 BIT IMAGES.jnl

Arithmetic	Interactive: Opens Arithmetic window for user to input background subtraction level from START IMAGE ANALYSIS.jnl step 3
Run Journal	Runs OPEN OBJECT LOG DDE FILE.JNL
Loop for all Images in a Directory	Interactive: Runs IMA OBJECTS 16 bit.jnl. User picks directory from which to choose.

IMA OBJECTS 16bit.jnl

Low Pass	3x3 convolution
Low Pass	3x3 convolution
Copy to 8-bit Image	No autoscale, to new untitled image
Save Using Sequential File Name	Saves 8bit image using previously defined Sequential File names.
Arithmetic	This background subtraction value needs to be manually entered into this journal from the value determined in START IMAGE ANALYSIS16 TO 8 BIT.jnl step 5
Threshold Image	1 unit above 0 to 255
Integrated Morphometry – Load State	Hoecsht.IMA Classifier 200 < area < 200000
Integrated Morphometry – Measure	Measures statistical info for all objects
Run Journal	Runs log obj and sum data.jnl

START IMAGE ANALYSIS 16 to 8 BIT.jnl

Copy to 8-bit Image	No autoscale, to new untitled image
---------------------	-------------------------------------

Close	Closes 16 bit image
Low Pass	3x3 convolution
Low Pass	3x3 convolution
Show Region Statistics	Interactive: Show entire image statistics. Calculate background subtraction value for step 6. by: INTENSITY Average + INTENSITY Std. Dev.
Arithmetic	Interactive: User inputs subtraction value from 5. into the constant Value field
Threshold image	Creates threshold by 1 unit above 0 to 255
Integrated Morphometry – Load State	Loads Start Image Analysis.ima Classifier 100 < area < 200000
Integrated Morphometry – Measure	Interactive: Shows area summary information about all objects. The average number is used as the Standard Area in 10.
Object Standards - Set Object Standards	Interactive: User inputs average area value from 9. into Standard Area box to be used by automated IMA for all images

IMA OBJECTS WITH NEW LOG FILE.jnl

Run Journal	OPEN OBJECT LOG DDE FILE.JNL
Run Journal	IMA OBJECTS.jnl
Close Summary Log	
Close Object Log	User must manually save every Excel spreadsheet generated.

INTERACTIVE IMA OBJECTS.jnl

Threshold Image	User manually sets threshold
Integrated Morphometry – Load State	Hoechst.IMA Classifier 200 < area < 200000
Integrated Morphometry – Measure	Objects
Integrated Morphometry – Log Data	Into open object.log file

COLLECT INTERACTIVE IMA DATA.jnl

Close Object Lo g	
Open Object Log	Interactive
Annotate Log File	Interactive: experimental information that will go into the first line of the object log file
Loop for all Images in Directory	Runs INTERACTIVE IMA OBJECTS.jnl

CHANGE FILTER, COLLECT IMAGE, SAVE SEQUENTIAL FILE  
NAME.jnl

Stage (Log Position)	
ADC-Focus	
Show Message and Wait	Interactive – user presses Enter when done focusing
ADC – Acquire from Digital Camera	Hoechst
Save Using Sequential File Name	
Close	Close open image

COLLECT HOECHST AND FITC.jnl

Run Journal	FOCUS, COLLECT IMAGE, SAVE SEQUENTIAL FILE NAME.JNL
Run Journal	CHANGE FILTER, COLLECT IMAGE, SAVE SEQUENTIAL FILE NAME.jnl

3X3 IMAGE COLLECTION HOECHST FITC.jnl

Stage (Scan)	COLLECT HOECHST AND FITC.jnl
--------------	---------------------------------

AUTOMATED 3X3 IMAGE COLLECTION HOECHST FITC.jnl

Set Up Sequential File Names	Interactive: user sets up prefix name and image storage directory
Open Data Log	Excel DDL files
Annotate Log File	Interactive: experimental information that will go into the first line of the log file of stage positions
Stage (Go to Origin)	Origin is set as the center of well A1
Stage (Move to Absolute Position)	Offset to upper left hand corner of well (1410, 1621)
Stage (Log Position)	
Stage (Scan Wells)	Interactive: user picks wells to scan: runs 3X3 IMAGE COLLECTION HOECHST FITC.jnl

AUTOMATED IMAGE COLLECTION.jnl

Set Up Sequential File Names	Interactive: user sets up prefix name and image storage directory
Open Data Log	Opens a DDE (Excel) File
Annotate Log File	Interactive: experimental information that will go into the first line of the log file of stage positions
Stage (Go to Origin)	Origin is set as the center of well A1
Stage (Log Position)	
Stage (Scan Wells)	Interactive: user picks wells to scan: runs FOCUS, COLLECT IMAGE, SAVE SEQUENTIAL FILE NAME.JNL. Well to well travel = (-9035, -9035)

STARTUP.jnl

Install and Configure Devices	Open Stage Meta Devices
Set Live Video Channel	
Preferences	<p><u>Measure Objects</u>: Draw failed classifier objects, Exclude objects that touch the edge of the image, Enable Elliptical Fourier Parameters, <b>turn off</b> Warn users when measurement data will be erased</p> <p><u>Image Saving</u>: Save Tiff/stk using LZW compression</p> <p><u>Image Windows</u>: Use transparent thresholds.</p>

Configure Default Paths	C:\Metamorph Data C:\Metamorph Data\Commmon Settings
Load Journal Taskbar	Common.JTB

### Nested Journals

#### Automated 3x3 Image Collection

*Loop* 3x3 image collection

*Loop* focus, collect image, save sequential file name

#### Automated 3x3 image collection Hoechst FITC

*Loop* 3x3 image collection Hoechst FITC

*loop* Collect Hoechst and FITC

focus, collect image, save sequential file name

change filter, collect image, save sequential file name

#### Automated image collection

*Loop* focus, collect image, save sequential file name

#### Collect automated IMA data in one Spreadsheet

Open object log DDE file

*Loop* IMA objects

Log obj and sum data

#### Collect automated IMA data in one spreadsheet 16 bit images

Open object log DDE file

*Loop* IMA objects 16 bit

Log obj and sum data

Although the above has been generally described in terms of a specific user interface and software code, other user interfaces and code can also be used. One of

ordinary skill in the art would recognize many other variations, alternatives, and modifications.

Fig. 6 is a simplified diagram 600 of a cleaning and dispensing system according to an embodiment of the present invention. This system 600 includes a variety of elements such as a dispensing head 609, which is coupled to a plurality of pipettes 601. The pipettes input and output fluids or solutions from plate 603. The plate has a plurality of sites, each of which can be used to input cells or a combination of cells and solution. The system also has elements to house solutions 605, which are used to manipulate cell samples in the plate. The dispensing head is supported through a support member 607, which is sufficiently rigid to allow for movement of the head. The dispenser is coupled to the present system in a mechanical and electrical manner, which provides for a fully integrated system for providing cell samples to the imaging system according to the present invention.

Fig. 7A illustrates a representative block flow diagram of simplified process steps of a method for determining properties of a manipulation based upon effects of the manipulation on one or more portions of one or more cells in a particular embodiment according to the present invention. This diagram is merely an illustration and should not limit the scope of the claims herein. One of ordinary skill in the art would recognize other variations, modifications, and alternatives. In step 700, one or more samples of cells can be provided. These cells can be live, dead, or fixed cells, or cell fractions. The cells also can be in one of many cell cycle stages, including G0, G1, S, G2 or M phase, M phase including the following cell cycle stages: interphase, prophase, prometaphase, metaphase, anaphase, and telophase.

Cell components tracked in presently preferable embodiments can include proteins, protein modifications, genetically manipulated proteins, exogenous proteins, enzymatic activities, nucleic acids, lipids, carbohydrates, organic and inorganic ion concentrations, sub-cellular structures, organelles, plasma membrane, adhesion complex, ion channels, ion pumps, integral membrane proteins, cell surface receptors, G-protein coupled receptors, tyrosine kinase receptors, nuclear membrane receptors, ECM binding complexes, endocytotic machinery, exocytotic machinery, lysosomes, peroxisomes, vacuoles, mitochondria, Golgi apparatus, cytoskeletal filament network, endoplasmic reticulum, nuclear membrane, proteosome apparatus, chromatin, nucleolus, cytoplasm, cytoplasmic signaling apparatus, microbe specializations and plant specializations.



The following table illustrates some markers and cell components commonly used by embodiments according to the present invention. Other markers can be used in various embodiments without departing from the scope of the invention.

<b>Cell component</b>	<b>Marker</b>	<b>Disease State</b>
Plasma membrane (including overall cell shape)	Carbocyanine dyes Phosphatidylserine Various lipids Glycoproteins	Apoptosis-Cancer Apoptosis-Neural degenerative Ds
Adhesion complexes	Cadherins Integrins Occludin Gap junction ERM proteins CAMs Catenins Desmosomes	Thrombosis Metastasis Wound healing Inflammatory Ds Dermatologic Ds
Ion Channels and Pumps	Na/K Atpase Calcium channels Serotonin reuptake pump CFTR	Cystic fibrosis Depression Congestive Heart Failure Epilepsy
G coupled receptors	$\beta$ adrenergic receptor Angiotensin receptor	Hypertension Heart Failure Angina
Tyrosine kinase receptors	PDGF receptor FGF receptor IGF receptor	Cancer Wound healing Angiogenesis Cerebrovascular Ds
ECM binding complexes	Dystroglycan Syndecan	Muscular Dystrophy

Endocytotic machinery	Clathrin Adaptor proteins COPs Presenilins Dynamin	Alzheimer's Ds
Exocytotic machinery	SNAREs Vesicles	Epilepsy Tetanus Systemic Inflammation Allergic Reactions
Lysosomes	Acid phosphatase Transferrin	Viral diseases
Peroxisomes/Vacuoles		Neural degenerative Ds
Mitochondria	Caspases Apoptosis inducing factor F1 ATPase Fluorescein Cyclo-oxygenase	Apoptosis Neural degenerative Ds Mitochondrial Cytopathies Inflammatory Ds
Golgi Apparatus	Lens Culinaris DiOC6 carbocyanine dye COPs	
Cytoskeletal Filament Networks	Microtubules Actin Intermediate Filaments Kinesin, dynein, myosin Microtubule associated proteins Actin binding proteins Rac/Rho Keratins	Cancer Neural degenerative Ds Inflammatory Ds Cardiovascular Ds Skin Ds

Endoplasmic Reticulum	SNARE PDI Ribosomes	Neural degenerative Ds
Nuclear Membrane	Lamins Nuclear Pore Complex	Cancer
Proteosome Apparatus	Ubiquityl transferases	Cancer
Chromatin	DNA Histone proteins Histone deacetylases Telomerases	Cancer Aging
Nucleolus	Phase markers	
Cytoplasm	Intermediary Metabolic Enzymes BRCA1	Cancer
Cytoplasmic Signaling Apparatus	Calcium Camp PKC pH	Cardiovascular Ds Migraine Apoptosis Cancer
Microbe Specializations	Flagella Cilia Cell Wall components: Chitin synthase	Infectious Ds
Plant specializations	Choloroplast Cell Wall components	Crop Protection

Then, in a step 702, one or more samples of the manipulation can be provided to the cells. Manipulations can comprise one or any combination of chemical, biological, mechanical, thermal, electromagnetic, gravitational, nuclear, or temporal factors, for example. For example, manipulations could include exposure to chemical compounds, including compounds of known biological activity such as therapeutics or drugs, or also compounds of unknown biological activity. Or exposure to biologics that may or may not be used as drugs such as hormones, growth factors, antibodies, or

extracellular matrix components. Or exposure to biologics such as infective materials such as viruses that may be naturally occurring viruses or viruses engineered to express exogenous genes at various levels. Bioengineered viruses are one example of manipulations via gene transfer. Other means of gene transfer are well known in the art and include but are not limited to electroporation, calcium phosphate precipitation, and lipid-based transfection. Manipulations could also include delivery of antisense polynucleotides by similar means as gene transfection. Other genetic manipulations include gene knock-outs or gene mutations. Manipulations also could include cell fusion. Physical manipulations could include exposing cells to shear stress under different rates of fluid flow, exposure of cells to different temperatures, exposure of cells to vacuum or positive pressure, or exposure of cells to sonication. Manipulations could also include applying centrifugal force. Manipulations could also include changes in gravitational force, including sub-gravitation (the preferred embodiment in outer space). Manipulations could include application of a constant or pulsed electrical current. Manipulations could also include irradiation. Manipulations could also include photobleaching which in some embodiments may include prior addition of a substance that would specifically mark areas to be photobleached by subsequent light exposure. In addition, these types of manipulations may be varied as to time of exposure, or cells could be subjected to multiple manipulations in various combinations and orders of addition. Of course, the type of manipulation used depends upon the application.

Then, in a step 704, one or more descriptors of a state in the portions of the cells in the presence of the manipulation can be determined using the images collected on the imaging system. Descriptors can comprise scalar or vector values, representing quantities such as area, perimeter, dimensions, intensity, gray level, aspect ratios, and the like. Other types of descriptors include, but are not limited to, one or any combination of characteristics such as a cell count, an area, a perimeter, a length, a breadth, a fiber length, a fiber breadth, a shape factor, a elliptical form factor, an inner radius, an outer radius, a mean radius, an equivalent radius, an equivalent sphere volume, an equivalent prolate volume, an equivalent oblate volume, an equivalent sphere surface area, an average intensity, a total intensity, and an optical density. These descriptors can be average or standard deviation values, or frequency statistics from the descriptors collected across a population of cells. These descriptors can be further reduced using other methods such as principal component analysis and the like. In some embodiments,

the descriptors include features from different cell portions or cell types. That is, a first feature can be from a nuclei and a second feature is from another cell structure such as Golgi apparatus, mitochondria, spacing between cell structures or cells themselves, as well as many others.

A presently preferable embodiment uses descriptors selected from the following table. Other descriptors can also be used without departing from the scope of the invention.

Name of Parameter	Explanation/Comments
Count	Number of objects
Area	
Perimeter	
Length	X axis
Width	Y axis
Shape Factor	Measure of roundness of an object
Height	Z axis
Radius	
Distribution of Brightness	
Radius of Dispersion	Measure of how dispersed the marker is from its centroid
Centroid location	x-y position of center of mass
Number of holes in closed objects	Derivatives of this measurement might include, for example, Euler number (= number of objects - number of holes)
Elliptical Fourier Analysis (EFA)	Multiple frequencies that describe the shape of a closed object
Wavelet Analysis	As in EFA, but using wavelet transform
Interobject Orientation	Polar Coordinate analysis of relative location
Distribution Interobject Distances	Including statistical characteristics
Spectral Output	Measures the wavelength spectrum of the reporter dye. Includes FRET
Optical density	Absorbance of light

Phase density	Phase shifting of light
Reflection interference	Measure of the distance of the cell membrane from the surface of the substrate
1,2 and 3 dimensional Fourier Analysis	Spatial frequency analysis of non closed objects
1,2 and 3 dimensional Wavelet Analysis	Spatial frequency analysis of non closed objects
Eccentricity	The eccentricity of the ellipse that has the same second moments as the region. A measure of object elongation.
Long axis/Short Axis Length	Another measure of object elongation.
Convex perimeter	Perimeter of the smallest convex polygon surrounding an object
Convex area	Area of the smallest convex polygon surrounding an object
Solidity	Ratio of polygon bounding box area to object area.
Extent	proportion of pixels in the bounding box that are also in the region
Granularity	
Pattern matching	Significance of similarity to reference pattern
Volume measurements	As above, but adding a z axis

Then, in a step 705, a database of cell information can be provided. Next, in a step 706, a plurality of descriptors can be searched from a database of cell information in order to locate descriptors based upon one of the descriptors of the manipulation. Then, in a step 708, properties of the manipulation are predicted based upon the properties of the located descriptors. Properties can comprise toxicity, specificity against a subset of tumors, mechanisms of chemical activity, mechanisms of biological activity, structure, adverse biological effects, biological pathways, clinical effects, cellular availability, pharmacological availability, pharmacodynamic properties, clinical uses and indications, pharmacological properties, such as absorption, excretion, distribution, metabolism and the like.

In a particular embodiment, step 706 comprises determining matching descriptors in the database corresponding to a prior administration of the manipulation to the descriptors of the present administration of the manipulation. In a particular embodiment according to the present invention, combinations of measurements of scalar values can provide predictive information. A database can be provided having one or more "cellular fingerprints" comprised of descriptors of cell-substance interactions of drugs having known mechanisms of action with cells. Such descriptors can be analyzed, classified, and compared using a plurality of techniques, such as statistical classification and clustering, heuristic classification techniques, a technique of creating "phylogenetic trees" based on various distance measures between descriptors from various drugs. In this embodiment, numeric values for the descriptors can be used by comparison techniques. A phylogenetic tree can be created that illustrates a statistical significance of the similarity between descriptors for the drugs in the database. Because the drugs used to build the initial database are of known mechanism, it can be determined whether a particular scalar value in a descriptor is statistically predictive. Finally, a compound descriptor with no known mechanism of action can be queried against the database and be statistically compared and classified among the drugs in the database that the compound most resembles.

In a particular embodiment, relationships between measured morphological properties of images and physiological conditions can be determined. Relationships can include, for example, treatment of different cell lines with chemical compounds, or comparing cells from a patient with control cells, and the like. In a presently preferable embodiment, comparisons can be performed on acquired image features. Some embodiments can comprise statistical and neural network - based approaches to perform comparisons of various features. The foregoing is provided as merely an example, and is not intended to limit the scope of the present invention. Other techniques can be included for different types of data.

In some embodiments, classification, clustering and other types of predictive data analysis can be performed on features extracted from cell images. In a presently preferable embodiment, statistical procedures for comparisons, classification and clustering are performed on data obtained from imaging cells.

Fragments of data preparation and pre-formatting (S language):

```
>tmp.frame <- Generic.Summary
```

```
> names1 <- paste("Cell.line.5", tmp.names, sep=".")
> by.compound.matrix <- as.matrix(arranged.by.compound)
```

Example of the code for principal component analysis (data preparation) using S language:

```
all.data.princomp <- menuPrincomp(data =
  by.compound.matrix, scores = T, cor = "Correlation",
  na.action = T, print.short = T, print.importance = T,
  print.loadings = T, cutoff.loadings = 0.1, plot.screplot
  = T, plot.loadings = T, plot.biplot = T,
  plot.biplot.choices = c(1,2), predict.p = F)
```

Example of clustering using a divisive hierarchical clustering algorithm:

```
> div.hier.2.manhattan.cluster$call
diana(x = tmp.sum.by.comp, diss = F, metric =
  "manhattan",
  stand = T, save.x = T, save.diss = T)
```

Another embodiment utilizes existing tools for biological sequence similarity searches, classification, and phylogenetic analysis. In a particular embodiment, numbers in a numerical descriptor can be substituted by one or more of nucleic acid or amino acid codes according to a one of several sets of rules. Once converted into a corresponding nucleotide or amino acid sequence representation, the fingerprints can be analyzed and compared using software and algorithms known in the art for genetic and peptide sequence comparisons, such as GCG, a product of Genetics Computer Group, with company headquarters in Madison WI. Select embodiments comprising such approaches enable the use of a broad array of sophisticated algorithms to compare, analyze, and cluster gene and protein sequences. Many programs performing this task are known to those of ordinary skill in the art, such as for example, the PHYLIP (PHYlogeny Interference Package) a package of programs for inferring phylogenies (evolutionary trees) described in (Feldenstein, J. 1996 *Methods Enzymol* 266:418-427 and Feldenstein, J. 1981 *J. Mol. Evol.* 17(6):368-376).



Embodiments can perform such analysis based upon factors such as numerical value, statistical properties, relationships with other values, and the like. Further details of a step of manipulation are noted more particular below.

Fig. 7B illustrates a representative block flow diagram of simplified process steps for determining one or more descriptors of a state in the portions of the cells in the presence of the manipulation of step 704 of Fig. 7A in a particular embodiment according to the present invention. This diagram is merely an illustration and should not limit the scope of the claims herein. One of ordinary skill in the art would recognize other variations, modifications, and alternatives. In a step 712, an image of a cell portion is obtained. In some embodiments, the cell portion is visualized with a fluorescently labeled marker that is specific for the portion or portions of interest. A cell portion can include, for example, one or more of the following: nuclei, Golgi apparatus, and other features. The cell portion may vary in select embodiments according to the invention. Then, in a step 714, a digitized representation of the image obtained in step 712 is determined. In some embodiments, steps 714 and step 712 can comprise a single step. These embodiments use a digital imaging means such as a digital camera, to obtain a digital image of the target directly. Next, in a step 716, the digital representation of the image is processed to obtain image features. Image features can include such quantities as area, perimeter, dimensions, intensity, aspect ratios, and the like. Then, in a step 718 descriptors can be determined from the image features. Descriptors can comprise scalar or vector quantities and can comprise the image features themselves, as well as composed features, such as shape factor derived by a relationship  $4\pi * \text{area} / \text{perimeter}$ , and the like. Descriptors can also comprise statistical quantities relating to feature characteristics across a population of cells, such as a standard deviation, and average, and the like.

In a preferred embodiment, cells can be placed onto a microscope, such as a Zeiss microscope, or its equivalent as known in the art. A starting point, named Site A01, is identified to the microscope. A plurality of exposure parameters can be optimized for automated image collection and analysis. The microscope can automatically move to a new well, automatically focus, collect one or more images, at one or more wavelengths, move to a next well, and repeat this process for all designated wells in a multiple well plate and for multiple plates. A file having a size and an intensity distribution measurement for each color and rank for each well can then be

created for the images acquired. Based on this information, a user or a computer can revisit sites of interest to collect more data, if desired, or to verify automated analysis. In a presently preferred embodiment, image automatic focus and acquisition can be done using computer software controlling the internal Z-motor of the microscope. Images are taken using a 10x, 20x, or 40x air long working distance objectives. Sometimes multiple images are collected per well. Image exposure times can be optimized for each fluorescent marker and cell line. The same exposure time can be used for each cell line and fluorescent marker to acquire data.

Fig. 7C illustrates a representative block flow diagram of simplified process steps for obtaining images of cell portions of step 712 of Fig. 7B in a particular embodiment according to the present invention. This diagram is merely an illustration and should not limit the scope of the claims herein. One of ordinary skill in the art would recognize other variations, modifications, and alternatives. The method is generally outlined by the steps below:

(1). In a step 720, a sample is provided to the imaging device. Samples can be provided in 96 well plates and the like. The sample may be loaded into a microscope, such as a Zeiss microscope or equivalent.

(2). In a step 722, a set of optical filters is selected to shine light of the appropriate wavelength to illuminate the first sample, which may be contained in a first well designated A01.

(3). In a step 724, an automatic focusing procedure is performed for the site. In a particular embodiment, the internal z-motor of the microscope which is attached to the objective nosepiece is used for automatic focusing of the microscope. In an alternative embodiments, the plate holding the samples is moved to perform automatic focusing of the microscope, or focusing can be performed by moving optical components attached to the microscope and the like.

(4). In a step 726, images are collected for the site. Images can be collected for every color at every site. Present embodiments can provide images for up to four colors. However, embodiments are contemplated that can provide more colors by using either a monochromator coupled with excitation filters which are on a filter wheel, or by digitally separating overlapping fluorophores. Those knowledgeable in the field will know that given calibration images of single fluorophores, a look-up table can be devised which will allow for the digital removal of fluorescence bleed-through of

fluorescence which may occur in optical channels other than the one for which that filter has been optimized in instances of using more than one fluorophore at once. Cell growth and density information is also collected. Cell density is determined by what percentage of the area being imaged is inhabited by cells. In some embodiments, imaging can be facilitated using one or more biosensors, molecules such as non-proteins, i.e., lipids and the like, that are luminescently tagged. However, some embodiments can also use fluorescence polarization and the like. Fluorescence polarization is a homogeneous fluorescence technology where the excited state of the molecule lasts much longer than in normal fluorescence, taking seconds to minutes to reach equilibrium, obliterating the need to wash away fluorescence markers that are not specifically bound to a marker. Further, embodiments can detect differences in spectral shifts of luminescent markers. Some fluorescence markers, such as Nile Red sold by Molecular Probes of Eugene, OR, will change its emission peak wavelength depending on its environment. One can detect these changes by monitoring the level of fluorescence at both wavelengths and reading out at ratio of the two.

(5). In a step 728, a determination is made whether more fields of view need to be taken for a particular color. If this is so, then processing continues at step 726 at a new site. Otherwise, processing continues with a decisional step 730. Images can now be taken by repeating step 726. In a preferred embodiment 4 to 9 images are collected at each site.

(5). In a step 730, a determination is made whether more optical configurations need to be taken in order to obtain images for all differently-marked cell portions the sample. If this is so, then in a step 732 a new optical configuration is determined. Images for the new optical configuration can now be taken by repeating steps 726 and 728.

(6). In a decisional step 734, after all optical configurations and images for fields of view in a sample have been obtained, a determination is made whether any further samples remain to be analyzed. If so, a new sample is brought into view and processing continues with step 720. Otherwise, image processing is complete. In a presently preferable embodiment, image data can be stored on a CD ROM using a CD ROM burner, such as CRW4416 made by Yamaha of Japan. However, other mass storage media can also be used.

Fig. 7D illustrates a representative block flow diagram of simplified process steps for processing digitized representations of step 716 of Fig. 7B in a particular embodiment according to the present invention. This diagram is merely an illustration and should not limit the scope of the claims herein. One of ordinary skill in the art would recognize other variations, modifications, and alternatives. The method is generally outlined by the steps below:

(1). In a step 740, a digitized image input is preprocessed. Preprocessing might include, but is not limited to, such operations as background subtraction, thresholding, smoothing, adoptive filtering, edge enhancements, contrast enhancements, histogram equalization. A particular combination of preprocessing steps can be applied to images in successive steps or in parallel to copies of the image.

A simplified example of a smoothing and background subtraction procedure in a MatLab language is presented in computer code below:

```
function Isubtracted = cmBackgrSubtrl(I,k)
% cmBackgrSubtrl(I,k) - simple flat background (=modal*k)
subtraction
% Y = cmBackgrSubtrl(I, k) - image Y is generated by
% subtraction (with saturation) of modal pixel value of I
multiplied by k
% DEFAULT - k=1
%
if (nargin == 1)
    k=1;
end
if (size(k)~=1)
    error('cmBackgrSubtrl: parameter k should be a number.
Exiting...');
end

%modpixnum = floor(size(I(:),1)/2);
%sortedval = sort( double(I(:)) );
%modpixel = sortedval(modpixnum);
```

```
modpixel = median(double(I(:)));
bg = k*modpixel;
```

```
Isubtracted = mmsubm( uint8(I), uint8(round(ones(
size(I))*k*modpixel )) );
```

An example of a procedure for thresholding in computer code (MatLab) is presented below:

```
function thresh = GetThreshByPerim1(I, M)
% GetThreshByPerim1(I) Finds optimal thresholding value for
image I
% N = GetThreshByPerim1(I) Finds thresholding value N for
image I
% N = GetThreshByPerim1(I, M) - tests threshold values up
to M
% DEFAULT M = maximum pixel value in I
% note that GetThreshByArea is significantly faster
% finds a threshold value that causes the maximal change in
the
% total perimeter of the objects (Russ ????)
% see Matlab_Auto_threshold1_1-23-99.doc for more details
% Note: works somewhat better on SMOOTH images (i.e.
medfilt2(I, [3 3]) two times

if (nargin == 0)
    error (strcat( mfilename, ' : at least one parameter
required')));
elseif (nargin == 1)
    M = double(max(I(:)));    %test thresholds up to
maximum pixel value in I
elseif (nargin > 2)
    error (strcat (mfilename, ' : too many parameters')));
end
```

```

if (size(M)>1)
    error (strcat(mfilename, ' : argument M should be a
number'));
end

Minval = double( min(I(:)));
step = 1;

%generate vertical vector perims with total perimeters of
objects at different
%threshold values
for i=Minval : step : M
    bwI = im2bw(I, i/255);
    prI = bwperim(bwI);
    pr = sum(prI(:));
    if (exist('perims', 'var') == 0) %perims is yet
undefined
        perims = pr;
    else
        perims = cat(1, perims, pr);
    end
end

% vector prdiffs contains differences between successive
perimeters
prdiffs = diff(perims);
mindecrease = min(prdiffs);
minvalues = find(prdiffs == mindecrease);
index_of_mindecrease = minvalues(1);
thresh = index_of_mindecrease + 1;

% =====end GetThresh1=====

```

Thresholding provides a specific intensity, such that pixels darker than the threshold are deemed black, and pixels lighter than the threshold are considered white. The thresholded image can be processed using binary image processing techniques in order to extract regions.

(2). In a step 742+744, the digitized image input is subjected to object identification. This can be accomplished by a variety of procedures, for example by thresholding or edge detection and subsequent morphological opening and closing. Edge detection can be accomplished by means of gradient-based or zero-crossing methods, such as Sobel, Canny, Laplacian, Perwitt, and other methods.

An example of object identification procedure based on Canny edge detection (in MatLab language) is presented below:

```
function Imask = cmMaskDNA1(I);
% cmMaskDNA1 - generates binary mask for cell nuclei
through edge detection
% Imask = cmMaskDNA1(I)
% PARAMETERS
%   I - intensity image (grayscale)
% OUTPUT
%   Imask - BW image with objects from I
%
% For more details see Notebook Matlab_DNA_masking1_1-22-
99.doc
% Uses SDC Morphology Toolbox V0.7

if (nargin ~= 1)
    error('Wrong number of input parameters');
end
if (nargout ~= 1)
    error('Wrong number of output parameters: one output
argument should be provided');
end
```

```

Imask = edge(I, 'canny');
Imask = mm dil(Imask, mmsecross(1));
Imask = mm ero ( mmclohole(Imask,mmsecross(1)));
Imask = mm edgeoff(Imask, mmsecross(1));
% note that mm edgeoff this command removed FILLED OBJECTS
but not touching OUTLINES.
% these outlines can be removed by filtering:
Imask = medfilt2(Imask, [5 5]);

%=====end cmMaskDNA1 =====

```

However, embodiments can also use other techniques, such as Fast Fourier Transforms (FFT) and the like as known in the art without departing from the scope of the present invention.

(3). In a step 746, a plurality of region features can be determined. For example, in a representative embodiment, image features can include such quantities as area, perimeter, dimensions, intensity, aspect ratios, and the like. Features not directly related to individual objects are also being extracted.

An example of a procedure for extraction of some of the features (MatLab language) is presented below:

```

function OData = cmGetObjectsData(I, Ilabel)
% cmGetObjectsData returns array measurements of objects in
image "I" masked by "Ilabel"
% EV 2-3-99; 2-10-99
% OData = cmGetObjectsData(I, Ilabel) returns an array of
morphological and intensity measurements
%   taken from a grayscale image "I". Objects are
identified on a mask image Ilabel, usually
%   created by bwlabel()
% OUTPUT:

```



```

% Each row in the output array OData represents individual
object
% columns contain the following measurements:
%
% 1 - Index ("number" of an object);      8 - Solidity;
% 2 - X coordinate of the center of mass; 9 - Extent;
% 3 - Y coordinate      -"-      ; 10 - Total
Intensity;
% 4 - Total Area (in pixels);              11 - Avg.
Intensity;
% 5 - Ratio of MajorAxis/MinorAxis;        12 - Median
Intensity;
% 6 - Eccentricity;                        13 - Intensity of
20% bright pixel
% 7 - EquivDiameter;                       14 - Intensity of
80% bright pixel
%
% For details on morphological parameters see information
on MatLab imfeature();
% Intensity parameters are either obvious or are documented
in comments in this file.

if (nargin ~= 2)
    error ('function requires exactly 2 parameters');
end
if (nargout ~= 1)
    error ('function has 1 output argument (array X by
14)');
end

% finished checking arguments

```

```
% first collect morphological parameters in a structure
array:
ImStats = imfeature(Ilabel, 'Area', 'Centroid',
'MajorAxisLength',...
'MinorAxisLength', 'Eccentricity', 'EquivDiameter', ...
'Solidity', 'Extent', 8 );

% now convert it into array (matrix) while collecting
intensity data for each object:

%preallocate output array:
numobjects = size(ImStats, 1);
OData = zeros(numobjects, 14);
%now convert ImStats into array and add intensity data to
it
for k=1:numobjects
    OData(k, 1) = k;
    OData(k, 2) = ImStats(k).Centroid(1);
    OData(k, 3) = ImStats(k).Centroid(2);
    OData(k, 4) = ImStats(k).Area;
    OData(k, 5) = (ImStats(k).MajorAxisLength) /
(ImStats(k).MinorAxisLength);
    OData(k, 6) = ImStats(k).Eccentricity ;
    OData(k, 7) = ImStats(k).EquivDiameter;
    OData(k, 8) = ImStats(k).Solidity;
    OData(k, 9) = ImStats(k).Extent;

    % now collect and assign intensity parameters from
image I

    object_pixels = find( Ilabel == k);
    object_area = size(object_pixels, 1); %same as total
number of pixels in the object
```

```

        object_intensities = double(I(object_pixels)); % need
to convert to double to do math
        sorted_intensities = sort(object_intensities); % will
need to get median, 20% and 80% pixels
        total_intensity = sum(object_intensities, 1);
        avg_intensity = total_intensity / object_area;
        median_intensity = sorted_intensities( floor(
object_area/2 ) + 1 );
        pix20 = sorted_intensities( floor(object_area*0.2)+1 )
; %brightest pixel among dimmest 20%
        pix80 = sorted_intensities( floor(object_area*0.8)+1 )
;

        OData(k, 10) = total_intensity;
        OData(k, 11) = avg_intensity;
        OData(k, 12) = median_intensity;
        OData(k, 13) = pix20; %brightest pixel among dimmest
20%
        OData(k, 14) = pix80; %dimmest pixel among brightest
20%
    end %for

%===== end function
cmGetObjectsData() =====

```

(4). In a step 748, quantitative descriptors, characterizing cell state are calculated based on the feature measurements extracted at step 746. For example, histogram distribution of intensities of cell nuclei provides information about the population cell cycle stages.

In a particular embodiment according to the present invention, data analysis techniques for describing the fluorescence patterns of cell portions in multiple cell lines in the presence and absence of compounds are provided. Automated image analysis techniques can include determining one or more regions from around nuclei,

individual cells, organelles, and the like, called "objects" using a thresholding function. Objects that reside on the edge of an image can be included or excluded in various embodiments. An average population information about an object can be determined and recorded into a database, which can comprise a database text file or Excel spreadsheet, for example. However, embodiments can use any recording means without departing from the scope of the present invention. Values measured can be compared to the visual image. One or more types of numerical descriptors can be generated from the values. For example, descriptors such as a number of objects, an average, a standard deviation of objects, a histogram (number or percentage of objects per bin, average, standard deviation), and the like can be determined.

In a particular embodiment according to the present invention, data can be analyzed using morphometric values derived from any of a plurality of techniques commonly known in the art. For example, a software package called MetaMorph Imaging System, provided by Universal Imaging Corporation, a company with headquarters in West Chester, PA and NIH Image, provided by Scion Corporation, a company with headquarters in Frederick, Maryland.

Fluorescent images can be described by numerical values, such as for example, an area, a fluorescence intensity, a population count, a radial dispersion, a perimeter, a length, and the like. Further, other values can be derived from such measurements. For example, a shape factor can be derived according to a relationship  $4\pi \cdot \text{area} / \text{perimeter}$ . Other values can be used in various embodiments according to the present invention. Such values can be analyzed as average values and frequency distributions from a population of individual cells.

In a particular embodiment according to the present invention, techniques for the automatic identification of mitotic cells are provided. Image analysis techniques employing techniques such as multidimensional representations, frequency-based representations, multidimensional cluster analysis techniques and the like can be included in various embodiments without departing from the scope of the present invention. Techniques for performing such analyses are known in the art and include those embodied in MatLab software, produced by MathWorks, a company with headquarters in Natick, MA.

Scalar values providing efficacious descriptors of cell images can be identified using the techniques of the present invention to perform predictive analysis of

drug behavior. In a presently preferred embodiment, a plurality of heterogeneous scalar values can be combined to provide descriptors for each manipulation. By applying predictive analysis routines to the collections of these descriptors, predictive information about any number of manipulations and cell interactions can be extracted.

Fig. 7E illustrates a representative block flow diagram of simplified process steps for analyzing image feature values to obtain descriptors of cell state of step 718 of Fig. 7B in a particular embodiment according to the present invention. This diagram is merely an illustration and should not limit the scope of the claims herein. One of ordinary skill in the art would recognize other variations, modifications, and alternatives. Fig. 7E illustrates an input data of descriptors of known manipulations 319. A step 320 of reformatting and transforming data 319 to formats suitable for analysis is performed. Additionally, a "cleaning" process can eliminate outlying data points and the like in the data. Then, in a step 322, a decision is made whether to continue with step 324 or with step 326 based upon determining a particular type of analysis appropriate for the present application or particular type of prediction. If decisional step 322 determines processing should continue with step 324, then, in that step, an error estimate using a set of test descriptors is performed to estimate the quality of a prediction and processing continues with step 320. Once an optimal prediction is achieved, processing continues with step 326. In step 326, optimal transformation parameters and prediction methods are selected for use in steps 328 and 330 which analyze data about an unknown manipulation. In a step 328, a solution is generated based upon any of techniques including training a neural network, solving a mathematical equation, applying decision tree rules and/or the like. In a step 330, an input data set of unknown descriptors 318 is reformatted and transformed based upon the optimal transformation parameters selected in step 326 using the transformation procedures in steps 320, 322 and 324. In a step 332, predictions techniques are applied to the reformatted manipulations from step 330 and the solution generated in step 328 and a plurality of properties of known manipulations 317 (e.g., therapeutic properties, and the like) in order to determine a prediction of properties of unknown manipulation 316.

Fig. 7F illustrates a representative block flow diagram of simplified process steps for a method of mapping a manipulation of cells to a physiological characteristic in a particular embodiment according to the present invention. This diagram is merely an illustration and should not limit the scope of the claims herein.

One of ordinary skill in the art would recognize other variations, modifications, and alternatives. The method is generally outlined by the steps below:

(1) In a step 750, a plurality of cells, e.g., dead, live, cell fractions or mixtures of cells are provided.

(2) Then, in a step 752, the plurality of cells is manipulated, where manipulation occurs using a source(s) from one or a combination selected from an electromagnetic, electrical, chemical, thermal, gravitational, nuclear, temporal, or a biological source.

(3) Next, in a step 754, a feature value is captured from the plurality of cells. The feature value can include one or any combination of characteristics such as cell count, area, perimeter, length, breadth, fiber length, fiber breadth, shape factor, elliptical form factor, inner radius, outer radius, mean radius, equivalent radius, equivalent sphere volume, equivalent prolate volume, equivalent oblate volume, equivalent sphere surface area, average intensity, total intensity, and optical density. This list is not meant to be limiting.

(4) Then, in a step 756, a degree of presence of one or more feature values is assigned for each manipulation.

(5) In a step 758, the feature values from the plurality of cells are stored in memory locations. From the memory locations the values can be used for statistical analyses to produce predictive information about the relatedness of the descriptors of the manipulations to one another. This information is used to infer properties of the manipulations.

Fig. 7G illustrates a representative block flow diagram of a simplified process steps for a method for populating a database with manipulated biological cell information in a particular embodiment according to the present invention. This diagram is merely an illustration and should not limit the scope of the claims herein. One of ordinary skill in the art would recognize other variations, modifications, and alternatives. The method is generally outlined by the steps below:

(1) In a step 760, a plurality of cells in various stages of the cell cycle, A montage image that was used as a source to generate data in Appendix A is presented in Fig. 12., such as for example, the stages of interphase, prophase, metaphase, anaphase, and telophase are provided.

(2) Then, in a step 762, each of the cells in the various stages of mitotic development is manipulated.

(3) Next, in a step 764, an image of the plurality of manipulated cells is captured using image acquisition techniques in order to provide a morphometric characteristic of each of the manipulated cells.

(4) As a preferable option, in a step 766, an image database may be populated with the image of the plurality of manipulated cells.

(5) Following step 764 or optional step 766, a morphological value is calculated from the image in a step 768.

(6) In a step 770, the database is populated with the morphological value.

Fig. 7H illustrates a representative block flow diagram of simplified process steps for a method for populating a database with manipulated biological information, e.g., image acquisition parameters, image feature summary information, and well experimental parameters in a particular embodiment according to the present invention. This diagram is merely an illustration and should not limit the scope of the claims herein. One of ordinary skill in the art would recognize other variations, modifications, and alternatives. Fig. 7H illustrates a step 780 in which cells are placed into site on a plate and a manipulation is applied. Then, in a step 781 an image is taken of the cells. In step 782, the image is transferred to an image archive database. Then, in a step 783, well experimental parameters are entered into the database 787. Well experimental parameters can include cell type, manipulation and the like. In a step 784, image acquisition parameters are transferred to database 787. Image acquisition parameters can include file name, fluorophores and the like. In a step 785, the image acquired in step 781 is analyzed. Then, in step 786, an image feature summary from the analysis step 785 is transferred to database 787.

In step 788, a lookup table for all analyses is provided to database 787. The lookup table provides information about the analyses. In a step 789, a query of database 787 for process data is performed. The results are reformatted. Then in a step 790, the database 787 is queried. Next, in a step 791, features of the manipulations stored in the database are combined and reduced. Next, in a step 793, reduced features of step 791 can be compared. In a step 792, the results of step 793 are recorded in database 787. Then, in a step 794, a report of predictions based on comparisons performed in step 793 is generated.

Fig. 7I illustrates a representative block flow diagram of simplified process steps for acquiring images of manipulated biological information, e.g., cells, cell tissues, and cell substituents in a particular embodiment according to the present invention. This diagram is merely an illustration and should not limit the scope of the claims herein. One of ordinary skill in the art would recognize other variations, modifications, and alternatives. Fig. 7I illustrates a step 770 in which a user sets up an image analysis procedure. Then, in a step 772, an image is read into image analysis software. Next, in a step 774, patterns and objects are identified in the image using one or more algorithms. Next, in a step 776, sets of features are extracted from the image. Then, in a step 778, feature information, descriptor values and the like are exported to the database, such as database 787 of Fig. 7H, for recording. Next, in a decisional step 779, a determination is made whether any more images should be taken. If this is so, processing continues with step 772. Otherwise, image acquisition processing is completed.

Fig. 7J illustrates a representative block flow diagram of simplified process steps for populating, acquiring and analyzing images of manipulated biological information in a particular embodiment according to the present invention. This diagram is merely an illustration and should not limit the scope of the claims herein. One of ordinary skill in the art would recognize other variations, modifications, and alternatives. Fig. 7J illustrates a step 300 of placing a plate onto an imaging stage and reading a bar code. Then, in a step 301 an autofocus procedure is performed. Next, in a step 302, a first optical filter configuration is selected and an image is collected. Then, in a decisional step 303, a determination is made whether more than one image per optical configuration can be taken. If so, then, in a step 304, a new position within the well is targeted and another image is collected. Then, in a decisional step 305, a determination is made whether any more images need to be collected. If this is so, step 304 is repeated until all images for a particular well have been collected. After one or more images are collected for the well, in a step 306, the stage is returned to a starting position within the well, and a montage is created from collected images. The results are named with a unique file name and stored.

In a decisional step 307, a determination is made whether any more optical channels in the well can be imaged. If this is so, then in a step 308 the next optical filter configuration is selected and an image is collected. Processing then



continues with decisional step 303, as described above. Otherwise, if no further optical channels in the well can be imaged, then in a decisional step 309 a determination is made whether any wells remain to be imaged. If not all wells have been imaged, then in a step 310, the stage moves to the next well and processing continues with step 301, as described above. Otherwise, if all wells on the plate have been imaged, then in a decisional step 311, a determination is made whether any more plates can be processed. If this is so, then processing continues with step 300 as described above. Otherwise, in a step 312, the information is stored to a CD or other storage device as a backup.

Fig. 7K illustrates a representative block flow diagram of simplified process steps compound based upon information about effects of one or more known compounds on a cell population in a particular embodiment according to the present invention. This diagram is merely an illustration and should not limit the scope of the claims herein. One of ordinary skill in the art would recognize other variations, modifications, and alternatives. Fig. 7K illustrates a step 340 of populating a database with descriptors for known compounds. Such descriptors can be determined from imaging the cell population. However, in some embodiments, descriptors can be derived by measurements and combinations of measurements and the like. Then, in a step 342, descriptors for the unknown compound are determined from imaging a second cell population. The second cell population has been treated with the unknown compound. Then, in a step 344, a relationship between the descriptors determined from the unknown compound with the descriptors determined from the known compounds can be determined. Finally, in a step 346, an inference can be made about the unknown compound based upon the descriptors of the known compounds from the relationship determined in step 344.

Accordingly, the present invention provides a novel database design. In a particular embodiment according to the present invention, a method for providing a database comprises measurement of a potentially large number of features of one or more sub-cellular morphometric markers. Markers can be from any of a large variety of normal and transformed cell lines from sources such as for example, human beings, fungi, or other species. The markers can be chosen to cover many areas of cell biology, such as, for example markers comprising the cytoskeleton of a cell. The cytoskeleton is one of a plurality of components that determine a cell's architecture, or "cytoarchitecture". A cytoarchitecture comprises structures that can mediate most

cellular processes, such as cell growth and division, for example. Because the cytoskeleton is a dynamic structure, it provides a constant indication of the processes occurring within the cell. The cytoarchitecture of a cell can be quantified to produce a one or more scalar values corresponding to many possible cellular markers, such as cytoskeleton, organelles, signaling molecules, adhesion molecules and the like. Such quantification can be performed in the presence and absence of drugs, peptides, proteins, anti-sense oligonucleotides, antibodies, genetic alterations and the like. Scalar values obtained from such quantification can provide information about the shape and metabolic state of the cell.

In a presently preferred embodiment, scalar values can comprise morphometric, frequency, multi-dimensional parameters and the like, extracted from one or more fluorescence images taken from a number of cellular markers from a population of cells. Two or more such scalar values extracted from a plurality of cell lines and markers grown in the same condition together comprise a unique "fingerprint" or descriptor that can be incorporated into a database. Such cellular descriptors will change in the presence of drugs, peptides, proteins, antisense oligonucleotides, antibodies or genetic alterations. Such changes can be sufficiently unique to permit a correlation to be drawn between similar descriptors. Such correlations can predict similar properties or characteristics with regard to mechanism of action, toxicity, animal model effectiveness, clinical trial effectiveness, patient responses and the like. In a presently preferred embodiment, a database can be built from a plurality of such descriptors from different cell lines, cellular markers, and compounds having known mechanisms of action (or structure, or gene response, or toxicity).

The present invention also provides database and descriptor comparisons according to other embodiments. In a particular embodiment according to the present invention, measurement of scalar values or features can provide predictive information. A database can be provided having one or more "cellular fingerprints" comprised of descriptors of cell substance interactions of drugs having known mechanisms of action with cells. Such descriptors can be compared using a plurality of techniques, such as a technique of creating "phylogenetic trees" of a statistical similarity between the descriptors from various drugs. In a present embodiment, scalar, numeric values can be converted into a nucleotide or amino acid letter. Once converted into a corresponding nucleotide representation, the descriptors can be analyzed and compared using software

and algorithms known in the art for genetic and peptide sequence comparisons, such as GCG, a product of Genetics Computer Group, with company headquarters in Madison WI. In an alternative embodiment, numeric values for the fingerprints can be used by comparison techniques. A phylogenetic tree can be created that illustrates a statistical significance of the similarity between descriptors for the drugs in the database. Because the drugs used to build the initial database are of known mechanism, it can be determined whether a particular scalar value in a descriptor is statistically predictive. Finally, a compound fingerprint with no known mechanism of action can be queried against the database and be statistically compared and classified among the drugs in the database that the compound most resembles.

In a particular embodiment, relationships between measured morphometric properties and features of images and physiological conditions can be determined. Relationships can include, for example, treatment of different cell lines with chemical compounds, or comparing cells from a patient with control cells, and the like. In a presently preferable embodiment, a clustering can be performed on acquired image descriptors. Some embodiments can comprise statistical and neural network - based approaches to perform clustering and comparisons of various descriptors. The foregoing is provided as merely an example, and is not intended to limit the scope of the present invention. Other techniques can be included for different types of data. In some embodiments, clustering and comparing can be performed on features extracted from cell images. In a presently preferable embodiment, procedures for comparisons and phylogenetic analysis of biological sequences can be applied to data obtained from imaging cells.

Select embodiments comprising such approaches enable the use of a broad array of sophisticated algorithms to compare, analyze, and cluster gene and protein sequences. Many programs performing this task are known to those of ordinary skill in the art, such as for example, the program Phylip, available at <http://evolution.genetics.washington.edu/phylip.html>, and other packages listed at <http://evolution.genetics.washington.edu/phylip/software.html>. However, select embodiments according to the present invention can comprise a technique of statistical classification, statistical clustering, distance based clustering, linear and non-linear regression analysis, self-organizing networks, and rule-based classification.

Embodiments can perform such analysis based upon factors such as numerical value, statistical properties, relationships with other values, and the like. In a particular embodiment, numbers in a numerical descriptor can be substituted by one or more of nucleic acid or amino acid codes. Resulting "pseudo-sequences" can be subjected to analysis by a sequence comparison and clustering program.

Other types of databases can also be provided according to other embodiments. The database includes details about the properties of a plurality of standard drugs. When the descriptor of a test compound is compared to the database, predictions about the properties of the test compound can be made using any known property of the other compounds in the database. For example, properties about a compound in the database could include structure, mechanism of action, clinical side effects, toxicity, specificity, gene expression, affinity, pharmacokinetics, and the like. The descriptor of a compound of unknown structure from a natural products library could be compared to the descriptors of compounds with known structure and the structure could be deduced from such a comparison. Similarly, such information could lead to better approaches to drug discovery research including target validation and compound analogizing, as well as pre-clinical animal modeling, clinical trial design, side effects, dose escalation, patient population and the like.

According to the present invention, databases can be integrated with and complementary to existing genomic databases. Differential genomic expression strategies can be used for drug discovery using database technology. In one particular embodiment, cell data and cellular response data can be associated with a genetic expression profile assay to form a single assay. Live cells expressing fluorescence markers can be treated with a drug, imaged and analyzed for morphometry; and then analyzed for mRNA for expression. Such embodiments can provide rapid development of tools to link cellular behavior with functional genomics.

Database methods according to the present invention can be used to predict gene function and to assist in target validation. Databases that include genetic diversity, i.e., having cellular descriptors from cells of differing genetic backgrounds (tumor, tissue specific, and gene knock out cell lines), can provide the capability to compare cells of unknown genetic background to those in the database. Similarly, the descriptor of an unknown cellular portion in the presence of multiple drugs can be queried against the descriptors of the known markers in the database. For example, if an

unknown gene is tagged with Green Fluorescent Protein (GFP), the database may be used to identify the cellular portions for which that unknown gene encodes.

According to the present invention, target validation and specialized cell-based assay screening can be performed using database systems and methods to serve as a universal high-throughput cell-based assay that can evaluate the molecular mechanism of drug action. As new genes are isolated and identified, a large collection of available gene-based knowledge is becoming available. From this large collection of new genes, potential protein targets can be identified using the genomic tools of sequence analysis and expression profiling. However, unless a gene mutation is tightly linked to a disease state, further validation of individual targets is a time consuming process, becoming a bottleneck in drug discovery. Furthermore, robotics and miniaturization are making "High Throughput Screening (HTS)" the industry standard, substantially reducing the time and cost of running a target-based biochemical assay. Therefore, it is now possible to routinely screen large libraries and use a resulting "hit" to validate the target. In such approaches, a specialized cell-based assay would be developed to test hits for each target. Since this often involves the creation of cell lines expressing new markers, this stage may also become a bottleneck that cannot keep pace with HTS. In addition, these cell-based assays may not be amenable to high-throughput screening, making it difficult to test the increasing number of analogs arising from combinatorial chemistry.

In a particular embodiment according to the invention, a rapid characterization of large compound libraries for potential use as pharmaceutical products can be provided by predicting properties of compounds that relate to the compounds' potential as bioactive drugs. In many drug discovery situations, virtually millions of compounds can be passed through a HTS assay against a small number of validated targets. These assays produce hundreds to thousands of potential hits. These hits can then be subsequently screened by a pipeline of secondary and tertiary screens to further characterize their specificity, often time completely missing non-specific interactions with other proteins. Techniques according to the present invention can provide a replacement to such screening operations by providing information about cellular accessibility and mechanism of action for the hits coming from a HTS system. Furthermore, it can replace the biochemical HTS assay and allow rapid and accurate identification of attractive compounds from large libraries without an intervening

biochemical assay. The cell information can be predictive of whether to continue into an animal model for each compound, and which animal model to pursue.

The principles of the present specifically contemplate a wide variety of research methodologies, or usage scenarios, implementing these principles. The following discussion of three such scenarios is by way of illustration and not limitation. Study of the principles enumerated herein will render evident to those skilled in the art certain additional methodologies or usage scenarios enabled by the teachings hereof. The present invention specifically contemplates all such modifications. The following description presents some specific embodiments and scenarios that represent a broader use of cellular phenotypic data and characterizations to deduce mechanisms of action and other features of cellular responses to various stimuli. Such procedures generally involve producing a quantitative cellular phenotype based upon two or more cellular attributes and then comparing that phenotype to phenotypes previously stored and indexed. Such procedures make use of databases or other repositories of biological information. The invention is not limited to the specific embodiments described here.

Considering first the procedure 2000 depicted in Figure 20, a compound has been identified as having a particular cellular activity. See 2004. For example, a compound may be found to inhibit the growth of certain cancer cell *in vitro* by a specific and desired mechanism of action. This may be a particular company's "gold standard."

Next, the compound is analyzed at 2006 in terms of its effect on one or more cell lines. More specifically, the compound is linked, virtually, to a particular phenotype. Two or more values or measures of cellular attributes characterize that phenotype. These attributes are quantified in the context of specific cellular markers.

In one example, the cellular marker is an organelle such as a nucleus or Golgi apparatus. Measured attributes useful for characterizing an associated phenotype include geometric parameters (e.g., size, shape, and/or location of the organelle) and composition (e.g., concentration of particular biomolecules within the organelle).

The phenotype may be characterized by administering the compound of interest to various cell lines and in various concentrations. In each example within this matrix, the attributes of interest are measured. Ultimately, certain phenotypic features (combinations of attribute values) are associated with the compound of interest. These features provide a template for the phenotype.

Next, using the phenotype as identified at 2006, the process identifies other compounds providing similar features. The goal here is to present a list of compounds having a mechanism of action similar to that of the compound that started the process. This allows researchers to identify a mechanism of action, if not already known, for their compound and to draw conclusions based upon their compound's link to other known compounds (which may not be chemically/structurally similar to the compound of interest).

Identifying similar compounds based upon phenotype can take many paths. Most will involve some mathematical basis. For example, the phenotype defined at 2006 can be represented as a fingerprint or vector comprised of multiple scalar values of cellular attributes (as described above). The phenotype representation can then be compared against known phenotypes characterized by the same format (e.g., they are all characterized as vectors having the same attribute set, but with different values of the attributes). The comparison may be as simple as a Euclidean distance or more sophisticated as a neural network or multivariate statistical correlation.

The known compounds and associated phenotypes may be stored as database records or other data structures that can be queried or otherwise accessed as part of the identification procedure. The compounds may also be associated with other relevant data such as clinical toxicity, cellular toxicity, hypersensitivity, mechanism of action, etc. (when available).

Compounds found to be sufficiently similar to the starting compound are returned for consideration by researchers. A data processing system may rank such compounds based on degree of similarity to the starting compound. In some cases, the system may even provide similarity scores associated with the listed compounds.

Often researchers wish to determine whether their particular compound has clinical or biochemical effects beyond those that they are already aware of. In a typical scenario, the compound of interest was selected based upon its strong binding a target or its stimulation or inhibition of cell growth in a particular cell line. The process associated with 2010 has likely identified the compound of interest as having a particular mechanism of action based on phenotypic similarity to other compounds having a similar mechanism of action. However, within the region of biochemical space, there may be subspaces (characterized by subphenotypes) that correspond to separate properties. For example, within the phenotypic space associated with one mechanism of action, there

may be subspaces associated with clinical toxicity, cellular toxicity (likely overlapping the clinical toxicity space), and little or no toxicity. Obviously, a researcher would like to know whether her compound is likely to be toxic.

Thus, the process 2000 may include characterizing the compound of interest in terms of its distance from (i.e., similarity to) specific phenotypes having known characteristics. In a typical example, the known characteristic is toxicity. This feature allows the researcher to quantify her compound in terms of mechanism of action AND toxicity (or in terms of two or more other relevant properties associated with phenotype). To allow simple ranking or characterization, compounds of interest may be scored according to a simple or weighted Boolean expression.

A second scenario of interest is depicted in Figure 21. This scenario again defines a phenotype in terms of a quantifiable vector or other measure. However, rather than using a compound of interest to generate the phenotype, some other cellular stimulus is used to generate the phenotype.

As shown, a process 2100 begins with receipt of cells of interest. See 2104. In many situations, the cells are produced by a genetic or epigenetic process that affects the expression level or activity of a particular protein. More generally, any cellular stimulus (e.g., radiation level and type, gravity level, magnetic field, acoustic perturbations, etc.) can be used to generate the cell line of interest. Importantly, this stimulus affects the phenotype and can be correlated therewith.

In the context of drug discovery, a gene encoding for a particular target can be genetically knocked out, underexpressed, overexpressed, expressed in a non-native state, etc. This may be accomplished via standard procedures involving genomic modification, translation or transcription apparatus modification (e.g., use of antisense nucleic acids), blocking target activity (using antibodies to a receptor site for example), and the like. These processes will generally affect the phenotype in some quantifiable way. Importantly, they clearly and unambiguously define a cellular phenotype associated with altering the activity of the target protein.

At 2106, the process involves measuring one or more cellular features from the cell line of interest to define/quantify the phenotype. This may be accomplished as described above with reference to 2006. Next, at 2108, the cellular phenotype generated in this manner is used to identify and rank a set of compounds



associated with the phenotype. This operation may proceed in the manner of operations 2008 and/or 2010 from Figure 20.

Finally, at 2110, the process clusters the compounds returned at 2108 by a mechanism of action. The operation 2106 has tightly bound a mechanism of action to a phenotype. Various compounds characterized and stored in a system database may be tentatively assigned a mechanism of action or may have no suggested mechanism of action. By matching their virtual phenotype to the phenotype generated at 2106, one can create or strengthen an association between the compounds and mechanism of action relevant to the stimulus at 2104.

Considering now Figure 22, a third scenario is depicted. This scenario again involves using a virtual phenotype to glean information relevant to a mechanism of action or other cellular activity. In this case, assay data from a group of compounds (e.g., a primary or focused library) is used to elucidate a phenotype.

As shown, a process 2200 begins by identifying a target protein. See 2204. Then, at 2206, the process involves identifying positive and negative biochemical hits. More generally, this may involve ranking a number of compounds based upon their interaction with the target. In a specific case, the compounds are ranked based upon their binding affinities to or ability to inhibit the enzymatic activity of the target protein.

After the compounds have been characterized in some manner based upon their interaction with the target, they are used to define a cellular phenotype. See 2208. Generally, the techniques to accomplish are the same as described with reference to operation 2006 of Figure 20. In this case however, one may obtain a strong correlation between mechanism of action (involving the target) and phenotype by using multiple of the compounds identified at 2206. For example, some of the "best hits" may be administered to cell lines in various concentrations. And some of the least effective compounds may also be administered. Cellular attributes that are more strongly exhibited with increasing concentration of the best hits (and not exhibited or exhibited only weakly upon administration of the negative hits) can be used to define the virtual phenotype. In a related approach, compounds having widely varying levels interaction with the target are administered to cells. Those cellular attributes that vary linearly or at least monotonically with the degree of interaction between the target and compound represent attributes that can be used to define the virtual phenotype.

After the cellular phenotype has been defined, previously characterized compounds may be clustered with that phenotype. See 2210. As with operation 2110 of Figure 2, this may create or strengthen an association between a mechanism of action and various compounds in a database.

Finally, and optionally, procedure 2200 may provide a "higher resolution" mechanism of action for the compounds identified at 2206. See 2212. Presumably interaction with the target suggests a specific mechanism of action or at least some aspect of a mechanism of action. However, a given target may participate in a larger cellular mechanism of action – unknown to researchers. Further, a compound may that binds with the target may participate in multiple mechanisms of action – some of which do not involve the target. By linking the target (and its positive hits) to a particular phenotype, some of these additional cellular level activities can be elucidated. The defined phenotype may have been previously identified as associated with other mechanisms of action or higher resolution mechanisms of action. Thus, the phenotype identified at 2208 can be leveraged to generate a higher resolution mechanism of action at 2212.

As suggested in the above discussion, compounds and associated phenotypes may be stored as database records. Such databases can take on many flavors. In one example, a database includes various pieces of information relevant to oncology. Such database may include numerous compounds classified by cellular phenotype, mechanism of action, toxicity, etc. More specifically, the database may include data on commercially available compounds clustered by cellular phenotypes corresponding to mechanisms of action. Further the databases of interest may extended or combined (via standard relational tables and algebra for example) to include additional data such as pharmacology data, cellular genomics data, gene expression data, protein expression data, etc. In a specific example, the database includes measurements made on a subset of the NCI60 cell lines, using DNA, Golgi apparatus, and/or microtubules as markers for defining the phenotypes. Other data includes dosage response information, variation in effect over time, etc. The compounds populating the database could include known National Cancer Institute oncology study compounds. In a specific embodiment, the compound set includes some or all of the compounds mentioned in the article "A gene expression database for the molecular pharmacology of cancer," Nature Genetics, 24, pp. 236-244 (March 2000).

Various biological analyses may be conducted to develop additional information for characterizing compound mechanisms of action, etc. For example, a cell count analysis may be used to develop dose response curves, GI 50 data, etc. The cell cycle may also be analyzed to find out how various stages in the cycle vary in response to particular stimuli. The Golgi apparatus may be analyzed to determine whether it is in a normal state, a dispersed state, a diffused state, etc. As another example, tubulin may be analyzed to determine whether it is normal, de-polymerized, over-polymerized, bundled, etc. Obviously, combinations of such analyses may be performed. For example, properties of the Golgi apparatus or tubulin may be analyzed over one or more cell cycles.

In some embodiments, techniques according to the present invention can provide tools for the later stages of drug development such as clinical trial design and patient management. The properties of known drugs, such as clinical trial and patient response information, will be used in a similar fashion as the pre-clinical information to provide predictions about the properties of novel compounds. Because the human cell is the locus of drug action, a database containing drug-cell interactions will be able to provide predictive value for this aspect of drug development.

**As the above discussion indicates, a single marker can provide multiple pieces of biologically relevant information. Image analysis is particularly well suited to handle this sort of detailed information. The advantages of using image analysis in this context can be understood by considering how certain related technologies operate. Many of these technologies employ laboratory automation and digital imaging to perform many cellular assays. But all have their limitations.**

**Most importantly, conventional cellular assays (e.g., gene chips, plate readers, etc.) measure average values of a population of cells. Thus, a significant problem is that these technologies operate on composite data from collections of cells. Multivariate analysis on composites does not have the power of multivariate analysis on individuals. For example, if property A is in 20 percent of a population and property B is also in 20 percent of the population, it is still important to know if this is the same 20 percent, a different 20 percent, or overlapping 20 percents.**

**To further illustrate the shortcomings of some conventional processes, these processes will now be described in more detail.**

**Gene chips:** A treatment is applied to cells. The cells are processed to extract DNA, mRNA or RNA, the latter of which is reverse transcribed into cDNA and hybridized with the probe on the chip, and message levels are measured. By the nature of this process, it reports the average message level (of unique messages) of a population of cells. It is possible to imagine a gene chip profile for a single cell, but the process is not easily scalable to many separate cells.

**FACS (Fluorescein Activated Cell Sorter):** A treatment is applied to cells. The cells are stained and the fluorescent intensity of each cell is measured. This is a cell-by-cell process. It reports a single measurement for each fluorescent marker (1-4 markers in practice). The sensor is a PMT (photo-multiplier tube).

**HTS (High Throughput Screening):** A treatment is applied to cells. The cells are stained and a plate reader measures the fluorescent level of each well. This process reports average intensity values for the population of cells in each well – one measurement per fluorescent marker.

A fundamental distinction between each of these techniques and the present invention can be understood by considering three parameters:

**Measurements/Fluorescent Marker:**

**Measurements/Cell**

**Total Number of Measurements (per experiment)**

Note that an “experiment” is the application of a treatment to a line of cells- each replicate, or different treatment, or different cell line is a different experiment.

Note that a marker might take many different forms. For example, a marker can be a label built into the cellular genome (e.g. GFP-Green Fluorescent Protein), a cellular component itself having a marker property (e.g. Campothecin), a direct stain (e.g. Hoeschst), or a antibody stain, or something else. The key differentiating one marker from another is the emitted light frequency or other signal from a label.

**Gene Chips:**

**Measurements/Fluorescent Marker:** not relevant

**Measurements/Cell:** average of cell population

**Total Number of Measurements (per experiment): 1 per DNA or RNA sequence represented on gene chip (1000s)**

**FACS:**

**Measurements/Fluorescent Marker: 1**

**Measurements/Cell: 1 per marker**

**Total Number of Measurements (per experiment): number of markers \* number of cells**

**HTS:**

**Measurements/Fluorescent Marker: 1**

**Measurements/Cell: average of a population of cells (all cells in a well)**

**Total Number of Measurements (per experiment): number of markers (each well is an independent experiment)**

**Present invention:**

**Measurements/Fluorescent Marker: 3 or more**

**Measurements/Cell: number of markers \* (3 or more)**

**Total Number of Measurements (per experiment): number of cells \* number of markers \* (3 or more)**

Thus, there is far more information content in the present invention's cell-by-cell image analysis than in other current characterization technologies. The only other technology that considers information on a per cell basis (FACS) considers only a gross value (measured as total number of photons) for each marker. Image analysis allows one to do significantly more with a single marker.

Figure 23 generally depicts a process flow that describes certain general operations employed in this aspect of the present invention. As depicted, a process 2301 begins at 2302 where cells of interest are labeled with one or more agents (markers) that bind to the cell. Note that the markers are chosen bind to separate components of interest contained within the cell.

After one or more cells of interest have been appropriately labeled and prepared, they are imaged in a fashion that shows the location of marked cell

components. The imaging apparatus accomplishes this by detecting signals emitted from the markers.

At block 2304 in process 2301, a computational system obtains images for each of the one or more markers. Note that these images may be combined in a single digital representation that provides information (e.g., signal intensity) about each of the one or more markers at each pixel. Alternatively, the images may be provided as separate digital representations (separate images) for each marker.

After the images have been obtained at 2304, the computational system next uses the images to generate one or more descriptors on each of the one or more markers. This operation is depicted at block 2306 and is substantially similar to block 704 set forth elsewhere herein. Finally, at 2308, the system classifies one or more cells into a number of biologically relevant classes using the markers and associated image descriptors. The number of biologically relevant classes is preferably equal to the number of markers under consideration plus two. So if there are two markers under consideration, then the method preferably provides at least four biologically relevant classifications. The examples below will illustrate how interactions between markers are used to this advantage.

As indicated above, this aspect of the invention is particularly useful in characterizing the effect that a particular stimulus has upon one or more cells. Thus, it will often be necessary to expose the cells to a particular stimulus prior to imaging. Examples of interesting stimuli include exposure to a chemical agent, exposure to a biological agent, exposure to radiation, and combinations thereof, as listed above. The cell or cells may be exposed to such stimuli prior to, during, and/or after exposure to the labeling agents.

As indicated, the cells will be labeled with one or more markers. A first marker binds to a first cell component and emits a signal in proportion to the concentration of that first component. Similarly, a second marker binds to a second cell component and emits a signal in proportion to the concentration of that second cell component. Such markers will typically label all cells in a population of cells, such as those cells present in the well of an assay plate.

Cell components of interest in this aspect of the invention include just about any particular component of a cell. Such components may be specific biomolecules, portions of biomolecules, and/or organelles and other subcellular

structures. Many examples of these components are presented elsewhere herein. Examples of particularly interesting components include DNA, Golgi components, cytoskeletal proteins, and combinations of these. In a particularly preferred embodiment, the cells of interest are labeled for the following combination of cell components: DNA, Golgi and tubulin.

Any of a number of different types of descriptors may be made on the images of the markers. Most of these descriptors represent a statistical or morphological characterization of the marker within the cell. Some of these descriptors operate on the spatial distribution of the marker within the cell. Others rely on an intensity histogram for the marker in the associated image. Lists of appropriate descriptors and appropriate markers are set forth elsewhere herein.

To carefully characterize the marker, one preferably makes at least two biologically relevant measurements of that marker. Full characterization will sometimes require three or even more measurements. Some measurements rely on previously determined statistical distributions of descriptors or combinations of two descriptors. Others rely on pattern recognition of relationships among two or more descriptors. Still others rely on statistical distributions determined by associated control experiments. Numerous other approaches to measurement will be apparent to those of skill in the art.

One example involves using DNA as a marked component. An image of DNA within a cell can be used to provide at least the following information: (1) the number of cells in a population (each cell's nucleus appears as a discrete region), (2) the quantity of DNA in a cell and hence the cell's interphase state ( $G_1$  versus S versus  $G_2$ ), and (3) the condensation state of the DNA to allow discrimination between mitotic and interphase cells. Thus, the biologically relevant information obtained from markers may take the form of multiple distinct biological measurements.

Often valuable biological information can be found in the interactions between two or more cellular components. In one preferred embodiment, there are three marked components of interest: DNA, Golgi, tubulin and multiple image analysis results per marker.

Initially in an image analysis example, DNA can be used to identify cells (more specifically, the cells' nuclei). Next one can analyze Golgi close to the

nuclei identified initially to determine a characteristic of the Golgi. In a related approach, one starts with DNA and uses it to identify cells. Then the DNA is analyzed with another algorithm to identify mitotic cells. Next one can analyze tubulin in mitotic cells to determine a measure for mitotic spindles. In yet another example, one can start with DNA and use it to identify cells and then analyze the DNA with another algorithm to identify mitotic cells. The DNA can also be analyzed with another algorithm to identify G2 cells.

Although the above has generally been described in terms of specific hardware, software, and methods, it is understood that many alternatives can exist. In particular, the present invention is not limited to a particular kind of data about a cell, but can be applied to virtually any cellular data where an understanding about the workings of the cell is desired. Thus, in some embodiments, the techniques of the present invention could provide information about many different types or groups of cells, substances, and genetic processes of all kinds. Of course, one of ordinary skill in the art would recognize other variations, modifications, and alternatives. Some examples according to the present invention are provided below.

## **EXPERIMENTS**

To prove the principle and demonstrate the objects of the present invention, experiments have been performed to determine the effects of manipulations on cell structure using imaging and analysis techniques applied to a variety of situations. These experiments were performed by growing multiple cell lines in the presence of multiple compounds, or substances. Cells were fixed and stained with fluorescent antibodies or labels to multiple cellular portions. One or more images of the cells were then obtained using a digital camera. Descriptors were built by quantifying and/or qualifying patterns of one or more feature from each image in the cell lines under study. A database was built from the descriptors. As the database grows, it should be able to predict the mechanism of action of an unknown drug by comparing its effect with the effects of known compounds or to identify data clusters within large libraries of compounds.

In a first experiment, an automated method to count the number of cells and differentiate normal, mitotic, and apoptotic cells was created. Approximately, 5,000 HeLa cells were plated per well in a 96 well plate and grown for 3.5 days. The cells



were fixed with  $-20^{\circ}$  MEOH for 5 minutes, washed with TBS for 15 minutes, and then incubated in 5 mg/ml Hoechst 33342 in TBS for 15 minutes. Then, 72 images were collected with a 40x objective and 75 ms exposure time.

The analysis was performed on objects that met a certain size criteria that was based on 1) measuring the size of objects in the image that were clearly not cells and 2) excluding the first peak of the area histogram (Fig. 8B values 1-4654).

Histograms of the individual object data were generated for each type of feature. Fig. 8A shows the histogram for average intensity, and Fig. 8B shows histogram data for the area of each object. Fig. 8C shows the scatter plot of the average intensity vs. the area of all of the objects. The pattern of the scatter plot showed an interesting pattern: a large cluster of cells in one region of the graph, with a scattering of object points in other regions. Because mitotic structures are identified as particularly bright objects, most likely due to the biological fact that the chromatin is condensed, the original Hoechst images could be used to identify which cells were either undergoing mitosis, or otherwise looked abnormal. Manual inspection of 917 cells resulted in the classification of each object. Fig. 8D shows a graph where each type of cellular classification is delimited. This graph clearly shows that the mitotic nuclei are brighter than the interphase nuclei. Further, the different phases of the cell cycle can be separated using these two features. Figs. 8E-8F show bar graphs of the average and standard deviations of the areas and average intensities for each cell classification type. These graphs show that interphase nuclei are statistically less bright than mitotic nuclei and that telophase nuclei are statistically smaller than other mitotic nuclei.

Each image was thresholded to an intensity level of 20. A standard area value was set at 9500 pixels. Automated information gathering about all of the objects was done and collected into an Excel spreadsheet (for more information see, section on imaging system). The following information was recorded:

IMAGE NAME
OBJECT #
AREA
STANDARD AREA COUNT
PERIMETER

FIBER LENGTH
FIBER BREADTH
SHAPE FACTOR
ELL. FORM FACTOR
INNER RADIUS
OUTER RADIUS
MEAN RADIUS
AVERAGE INTENSITY
TOTAL INTENSITY
OPTICAL DENSITY
RADIAL DISPERSION
TEXTURE DIFFERENCE MOMENT
EFA HARMONIC 2, SEMI-MAJOR AXIS
EFA HARMONIC 2, SEMI-MINOR AXIS
EFA HARMONIC 2, SEMI-MAJOR AXIS
ANGLE
EFA HARMONIC 2, ELLIPSE AREA
EFA HARMONIC 2, AXIAL RATIO
EFA HARMONIC 3, SEMI-MINOR AXIS

The following results were obtained:

- 1,250 objects were counted
- 201 of those objects has standard area counts  $> 2$  (area  $> 19000$  pixels)
- 195 objects had areas  $< 6000$  pixels
- 1529 objects estimated in total
- 1328 object areas are  $> 6000$  pixels
- The data was reduced to 917 objects that were  $6000 < \text{area} < 19000$
- For the 917 objects a scatter plot of area vs. average intensity and a histogram of the average intensity were generated.
- 116 objects that had average intensity intensities  $> 60$  were manually looked at to determine their morphology.
- Of those 116 objects:

6 were dead or indistinguishable

4 were interphase

30 were prophase

32 were metaphase

24 were anaphase

20 were telophase (10 pairs)

- 12 prophase objects were missed because of gray scale cut off. (8 of those prophase cells had gray scale values  $> 57$ , as did 7 interphase)
- 1 telophase object was missed because it was too small ( $< 6000$ )
- 1 prophase object was missed because it was too big ( $> 1900$ )
- 16 mitotic objects were missed because they were parts of objects with standard count  $> 2$ .

In sum, out of 917 single objects, the analysis correctly identified 106 out of 130 mitotic objects, or (81% predictive, 91% of identified mitotics). Out of 917 single objects, the analysis incorrectly identified only 10 non-mitotics as mitotics (1% total, 8% of identified mitotics); 14 mitotics as interphase (1.4% total, 1% interphase). An automated classification system that would automatically assign values to each object using these or other measurement features can thus be developed, utilizing the principles set forth herein.

In a second experiment, the effects of Taxol on MDCK cells and the different types of morphological effects were observed. A plurality of MDCK cells grown in 96 well plates were treated with Taxol for 4.5 hours at different concentrations (10  $\mu$ M-1pM). They were then fixed, labeled with Hoechst, and imaged..

This experiment used a labeling protocol comprising: MEOH fix at  $-20^{\circ}$ , Wash in PBS, Block in PBS/BSA/Serum/Triton-X 100, Incubate with 5  $\mu$ g/ml Hoechst 10 minutes, and wash.

Cells were inspected for different morphologies and manually counted at each different drug concentration in one well. Fig. 9 shows example images from each drug concentration and the different types of morphologies and cells are highlighted. Fig. 10 shows the distribution of each morphology within the cell population as a function of drug concentration. The higher the concentration of Taxol, the larger

proportion of cells underwent apoptosis, and the fewer number of normal mitotic cells were detected.

In a third experiment, the purpose was to determine whether the automated analysis methods developed in the first experiment can detect differences in Hoechst morphology in the presence of 6 known compounds at one concentration and exposure time in one cell line. In this experiment, HeLa cells were separately treated with 6 compounds with known mechanism of action. The quantitative methods described in the first experiment were applied to the Hoechst images.

Approximately 5,000 HeLa cells per well were plated in a Costar black-walled 96 well tissue culture treated plate and left to recover in the incubator for 24 hours. After this time, 10 ug/mL of cytochalasin D (CD), Taxol, hydroxyurea, vinblastine, nocodazole, and staurosporine was added to different wells at a 1:100 addition in DMSO. The cells were incubated in the presence of drug for 24 more hours. After 24 hours, the cells were removed and fixed as in the first experiment. Then, 9 images per well were collected of the Hoechst staining using a 10x objective.

The low magnification images taken of Hoechst were run through the automated image analysis method described in the first experiment. Plots of the average intensity and area were made of each compound. Fig. 11 shows the scatter plots of the compounds. The scatter plots of each compound are visually distinct. For example, cells treated with CD are smaller than control, and cells treated with Hydroxyurea are larger and brighter. Furthermore, the number of cells per well was very different (data not shown).

The effects of different compounds can be clearly and automatically distinguished by identifying changes in cellular morphology. This method can also be used to count adherent cells.

The next experiment was to develop clustering algorithms that assign statistically meaningful values to the representative two dimensional data shown in Fig. 10, and even more complicated clustering of all of the multidimensional data that can be extracted across one, and multiple images.

A fourth experiment was performed to obtain high magnification images of two markers in the presence of drugs. In this experiment, HeLa cells were treated with 80 generic compounds with known mechanism of action. The quantitative methods described in the first experiment were applied to the Hoechst images.

Approximately 5,000 HeLa cells per well were plated in a Costar black walled 96 well tissue culture-treated plate and left to recover in the incubator for 24 hours. After this time, 10 ug/mL of each compound from the Killer Plate from Microsource Discovery Systems (Gaylordsville, CT) was added to different wells at a 1:100 addition in DMSO. The cells were incubated in the presence of drug for 24 more hours. After 24 hours, the cells were removed and fixed as in the first experiment. In addition to being labeled with Hoechst 33342 (against chromatin), cells were also labeled with 1 unit of rhodamine-conjugated phalloidin (against actin) for 30 minutes.

The 96 well plate was imaged twice. Once, 9 images per well were collected of the Hoechst staining using a 10x objective. After this, one image per well of both the phalloidin and Hoechst staining was collected using a 40x objective.

The resulting high magnification images were analyzed qualitatively and distinct pattern differences were detected in both the Hoechst and phalloidin images. Fig. 12 shows three example images from the experiment. The top row is the Hoechst staining, and the bottom row is the phalloidin staining from the same well. The columns show the images from wells treated with just DMSO (control), cytochalasin D, and Colchicine. The morphology of each marker is different in the presence of each drug. Interestingly, there is an effect in the morphology of the chromatin in the Hoechst image of cytochalasin D, which directly targets the actin cytoskeleton (and thus there is an expected effect in the phalloidin image). Also, there is an effect on the actin cytoskeleton, compared to control, in the presence of colchicine that directly targets the microtubule network.

The low magnification images were analyzed as described in the first experiment, and different patterns were seen in both the average intensity vs. area plots, and in the number of cells per well (data not shown). Thus, changes in patterns of a marker that is "down-stream" from the direct target of a compound are detectable. Automated image analysis protocols for actin and other markers can be developed similarly, again utilizing the principles set forth herein.

A fifth experiment was performed to test quadruple labeling of 9 different cell lines grown in normal conditions. In this experiment, NCI-H460, A549, MDA-MD-231, MCF-7, SK-OV-3, OVCAR-3, A498, U-2 OS, and HeLa cells were plated. Then, the cells were fixed and stained for portions of the each cell known as DNA, tubulin, actin, and Golgi.

The following table summarizes the procedures for this experiment:

Action	Active Ingredient/Notes	Buffer	Vol/ well	Desired Time	Temp
Remove media	NOTE: gently by pipetting, not aspiration				
Fix	4% Formaldehyde	PBS	100 $\mu$ l	20 min	rt
Wash		TBS	100 $\mu$ l	5 min	rt
Wash		TBS	100 $\mu$ l	5 min	rt
Permeablize	0.1% Triton X-100	TBS	100 $\mu$ l	10 min	rt
Permeablize	0.1% Triton X-100	TBS	100 $\mu$ l	10 min	rt
Block	% BSA % Serum Filter sterilize before use	TBS w/azide	100 $\mu$ l	1hr or o/n	rt or 4°C
Primary Antibody	1:1000 dilution of DM1 $\alpha$	TBS + 1% BSA + 0.1% TX-100	50 $\mu$ l	1hr or o/n	rt or 4°C
Wash		TBS	100 $\mu$ l	5 min	rt
Wash		TBS	100 $\mu$ l	5 min	rt
Wash		TBS	100 $\mu$ l	5 min	rt
Fluorescent Stain	FITC lens culinaris 1:500 Rhodamine-Phalloidin 1:500 CY5 goat anti-mouse 1:100	TBS + 1% BSA + 0.1% TX-100	50 $\mu$ l	1 hr.	rt, dark
Wash		PBS	100 $\mu$ l	5 min	rt, dark
Hoechst	1:1000 dilution of 5mg/ml	TBS	100 $\mu$ l	15 min	rt, dark
Wash		PBS	100 $\mu$ l	5 min	rt, dark

Wash		PBS	100 $\mu$ l	5 min	rt, dark
Wash		PBS	100 $\mu$ l	5 min	rt, dark
Store		PBS	200 $\mu$ l	1 month	4°C

Cells were plated out at different densities for 48 hours. Cells were fixed and labeled by the above method. Cells were imaged using an automated imaging system that collected 9 images from each marker using a 10x objective. Higher magnification images were collected of a few cells for demonstration purposes.

In this experiment, each cell line demonstrated different morphological patterns as determined by phase. For example, A549 cells are much more compacted than OVCAR-3 cells as determined by phase contrast imaging (data not shown). The different fluorescent markers showed even bigger differences between different cell lines. Figs. 13 and 14 show 4 panels of each marker for A549 (Fig. 13) and OVCAR-3 cells (Fig. 14). The markers are Hoechst (upper left), Phalloidin (upper right), Lens culinaris (lower left), and DM1 $\alpha$  antibody (lower right). The following table summarizes the qualitative differences between these images:

MARKER	A549	OVCAR3
Hoechst/DNA	small	large
Phalloidin/actin	fuzzy	crisp - many stress fibers
Lens culinaris/Golgi	compact	Disperse/punctate
DM1 $\alpha$ /Tubulin	perinuclear	evenly distributed

Higher magnification images were taken of the OVCAR3 cells. Fig. 15 shows the same markers at 20x, and Fig. 16 shows the markers at 40x. While the highest magnification images show the most detail, these images illustrate that very little morphological or feature information is lost in the 10x images.

These data exemplify the differences in morphology seen between different cell types. Thus the automated image analysis software can be customized for

each marker in each cell type. Different drugs should effect these morphologies differentially.

An automated quantification method for each marker and cell line can be similarly developed.

A sixth experiment was conducted with a more sophisticated software package and to develop more flexible image recognition algorithms. In this experiment, prototype image features extraction was performed using MatLab programming language with image toolbox and SDC morphology toolboxes. Algorithms are being developed that will automatically identify objects on images and to measure various morphological and feature parameters of these objects. Many different features for each of the cellular markers were acquired.

An example of a MatLab program called "AnalyseDNA" that takes as an input an unlimited number of images, identifies individual objects in these images based on either their intensities, or based on edge-detection algorithms, and extracts a number of morphological and intensity characteristics of these objects. A copy of this program follows:

**Listing of the AnalyseDNA.m program and of some of the  
supporting subroutines**

```
function files_analysed = AnalyseDNA(filemask, outpath, nx,
ny, filter_range, dext, modifier, sfname)
% AnalyseDNA performs measurements on files of DNA images
% V1. EV 2-11-99; 2-15-99; 2-16-99
%
% files_analysed = AnalyseDNA(filemask, outpath, nx, ny,
filter_range, dext, modifier, sfname)
%
% PARAMETERS:
%   ALL PARAMETERS ARE OPTIONAL
%
%   FILEMASK - mask for file names to be analyzed
INCLUDING PATH(for example c:\images\*.tif)
```



```
%    DEFAULT '*.tif' (all *.tif files in the current
directory).
%
%    OUTPATH - path to a directory where all the output
files will be placed.
%    DEFAULT - output is saved in the same directory which
contains images
%
%    NX, NY - number of individual images in montage images
along X and Y axes (DEFAULT 1)
%
%    FILTER_RANGE - 3 col-wide array (or []). Specifies how
data is filtered when summary is calculated
%    this parameter internally is passed to GetDNADData and
then to GetSummaryData - see these
%    functions for details. For example: [2 2 Inf; 6 100
8000] will case all rows of data for which
%        values in column 2 are less than 2 and all rows
where values in column 6 are less than 100 or
%        more than 8000 to be excluded from all
calculations of a summary.
%    DEFAULT - [] (means do not filter, summarize all data)
%
%    DEXT - string. Extension for data files being saved.
%    DEFAULT 'dat';
%
%    MODIFIER - this modifier is 'SUMMARY', summary file is
created;
%        'SUMMARY ONLY' - only summary is generated, data
for individual files are not saved
%
%    sfname - string. File name of a summary file
%    DEFAULT 'summary[date].dat'
```

```
%  
% OUTPUT:  
%  
%   AnalyseDNA works on image files or montages. For each  
image file it creates a tab-delimits file of measured  
%   parameters of all the objects in the montage with the  
same base name as a montage file and extension specified  
%   by dext parameter (or .dat by default) and file  
'errors[date].err' - with the list of files that matched  
the  
%   filemask but could not be processed.  
%   If 'summary' or 'summary only' modifier is specified,  
it also creates a single file 'summary[date].dat' (or  
%   different extension, if specified by DEXT) which  
contains summary information for all analyzed files.  
%  
%   ALL OUTPUT FILES are saved in a directory specified by  
OUTPATH parameter  
%  
%   RETURNS *files_analysed* - number of files that have  
been successfully processed.  
%  
%   Column designations in the output files are described  
in GetDNAData  
%  
% FILE NAME CONVENTIONS  
%   AnalyseDNA attempts to identify a number for each file  
to identify the file in summary output.  
%   It does that by looking for the first space or  
underscore, followed by a number and then takes  
%   as many successive numbers as it can find. If it fails  
to identify a number it assigns a  
%   default which is -1
```

```
%  
%  
% SEE ALSO GetDNADData, GetSummaryData  
%  
% TO DO    improve error handling in opening and writing  
files (GLOBAL error_file ?)  
%          include procedures for writing text headers into  
the output files  
  
if nargin > 8  
    error ('Wrong number of input parameters');  
end  
if nargout > 1  
    error ('Wrong number of output parameters: only one  
allowed');  
end  
  
% set defaults  
need_summary = 0;  
summary_only = 0;  
use_default_outpath = 0;  
datestring = datestr(floor(now));  
if nargin == 7      % set default summary file name  
    sfname = ['summary' deblank(datestring)]; % extension  
will be appended later based on dext  
    if deblank(upper(modifier)) == 'SUMMARY'  
        need_summary = 1;  
    elseif deblank(upper(modifier)) == 'SUMMARY ONLY'  
        need_summary = 1;  
        summary_only = 1;  
    else  
        error (['Wrong parameter: unknown modifier '  
modifier]);
```

```
end
end

if nargin == 5
    % default data file extension
    set dext = 'dat';
end
if nargin == 4
    % default filter range
    filter_range = [];
end
if nargin == 3
    ny = 1; % default number of images in montage along Y
end
if nargin == 2
    nx = 1;
end
if nargin == 1
    use_default_outpath = 1;
end
if nargin == 0
    filemask = '*.tif'
end

% check parameters
if ( ~ischar(filemask) | ~ischar(dext) | ~ischar(sfname) )
    error('Wrong parameter type: filename, filepath,
dext and sfname should be strings');
end
if ( ( size(nx) ~= [1 1] ) | ( size(ny) ~= [1 1] ) )
    error('Wrong parameter type: nx and ny should be
scalars (1x1 arrays)');
end
```

```
if (~isempty(filter_range) & size(filter_range, 2) ~= 3)
    error ('Wrong parameter type: filter range should be []
or 3 - cols-wide array');
end
% end testing parameters

% Generate list of files to process

datapath = getpath(filemask);
if use_default_outpath == 1
    outpath = datapath;
end
if exist(outpath, 'dir') ~= 7
    error(['Path ' outpath, 'not found. Exiting..']);
elseif exist(datapath, 'dir') ~= 7
    error(['Path ' datapath, 'not found. Exiting..']);
end

sfname = makefullname(outpath, sfname, dext);
if need_summary == 1
    if exist(sfname, 'file')
        disp(['File ', sfname, 'already exists!']);
        input ('Press ^C to abort, Enter to delete and
continue');
        delete(sfname);
    end
end

flist = FileList(getfname(filemask), datapath);
numfiles = size(flist, 1); % total number of files to
process
disp(['About to process ', num2str(numfiles), ' files']);
```

```
%DEBUG - commented out "input" to run from Wrod
input('Press ^C to abort, Enter to continue');

% main loop where the job gets done:
error_file = makefullname(outpath, ['error' datestring
'.err']);
num_processed = 0;
num_error = 0;
for i = 1:numfiles
    % first generate file name for a data output file
    current_fullname = flist(i, :); % full name with path
    and extension
    current_datafile = makefullname(outpath,
makefname(getbasefname(current_fullname), dext) );

    %extract number from a filename
    fnumber = getfilenumber(current_fullname);

    % load an imagefile, record errors
    read_error = 0;
    try
        I = imread(current_fullname);
        %DEBUG
        disp(['Image file #', num2str(fnumber), '
loaded']);
    catch
        % record file-opening error in an error_file
        read_error = 1;
        num_error = num_error +1;
        msg = [current_fullname ': ' lasterr];
        add_error_msg(error_file, msg);
    end
end
```

```

% extract and write data to a file in outpath
if read_error ~=1
    if (need_summary == 0)
        %DEBUG
        disp(['Starting analysis of file #',
num2str(fnumber), '.']);
        current_data = GetDNAData(I, nx, ny, fnumber);
        %DEBUG
        disp(['Finished analysis of file #',
num2str(fnumber), '.']);
        %load current_data.mat 'current_data';
        write_data(current_data, current_datafile);
    else %summary needed
        %DEBUG
        [current_data, current_summary] = GetDNAData(I,
nx, ny, fnumber, filter_range);
        %load current_data.mat 'current_data';
        %load current_summary.mat 'current_summary';
        write_summary (current_summary, sfname);
        if summary_only ~= 1
            write_data(current_data, current_datafile);
        end
    end
end
end % of the main for loop
num_processed = numfiles - num_error;

%=====end function AnalyseDNA()
=====

%=====
=====

function result = add_error_msg(filename, msg)

```

```

% adds string MSG to an errorfile FILENAME
% returns 1 if success, 0 if failure

err_FID = fopen(filename, 'at');
if err_FID == -1
    warning(['Can not open error file ' filename]);
else
    fprintf(err_FID, '%s\n', msg);
    fclose(err_FID);
end
%=====end function add_error_masg()
=====

%=====
=====

function N = getfilenumber(fname)
% returns the first number extracted from a file name
(string) or -1 if fails to extract any number
numbers = NumbersFromString( getfname(fname) ); % vector of
all numbers encoded in the name

                                % (but not in the path, even if
present)
if isempty(numbers)
    N = (-1); % return -1 if no numbers found in the name
else
    N = numbers(1);
end

%===== end function getfilenumber()
=====

```



```

%=====
=====
function result = write_data(data_array, file_name)
% writes data in a data_array in a tab-delimited ascii
file.
% result is 0 if success and -1 if failure
% if file_name exists, overwrites it
result = -1;
try
    fid = fopen(file_name, 'wt');
    if fid ~= -1
        for k = 1:size(data_array, 1)
            fprintf(fid, '%g\t', data_array(k, :));
            fprintf (fid, '\n');
        end
        test = fclose(fid);
        result = -1;
    catch
        result = -1;
    end

%===== end function write_data()
=====

%=====
=====
function result = write_summary (s_vector, file_name)
% appends summary vector s_vector to a file_name (ASCII
tab-delimited file).
% if file_name does not exist, creates it.
% result is 0 if success and -1 if failure
%
result = -1;

```

```

try
    % debug
    fid = fopen(file_name, 'at');
    result = fprintf(fid, '%g\t', s_vector);
    result = fprintf(fid, '\n');
    result = fclose(fid);
    result = 0;
catch
    result = -1;
end

% ===== end function write_summary()
=====

function Data = GetObjectsData(I, Ilabel)
% GetObjectsData returns array measurements of objects in
image "I" masked by "Ilabel"
% EV 2-3-99; 2-10-99
% OData = GetObjectsData(I, Ilabel) returns an array of
morphological and intensity measurements
%   taken from a grayscale image "I". Objects are
identified on a mask image Ilabel, usually
%   created by bwlabel()
% OUTPUT:
% Each row in the output array OData represents individual
object
% columns contain the following measurements:
%
%   1 - Index ("number" of an object);      8 - Solidity;
%   2 - X coordinate of the center of mass; 9 - Extent;
%   3 - Y coordinate      -"-      ; 10 - Total
Intensity;

```

```
%      4 - Total Area (in pixels);                11 - Avg.
Intensity;
%      5 - Ratio of MajorAxis/MinorAxis;          12 - Median
Intensity;
%      6 - Eccentricity;                          13 - Intensity of
20% bright pixel
%      7 - EquivDiameter;                        14 - Intensity of
80% bright pixel
%
% For details on morphological parameters see information
on MatLab imfeature();
% Intensity parameters are either obvious or are documented
in comments in this file.
% Procedures in this file are documented in notebook file
"MATLAB Measuring Nuclei (1) 1-29-98.doc"

if (nargin ~= 2)
    error ('function requires exactly 2 parameters');
end
if (nargout ~= 1)
    error ('function has 1 output argument (array X by
14)');
end

% finished checking arguments

% first collect morphological parameters in a structure
array:
ImStats = imfeature(Ilabel, 'Area', 'Centroid',
'MajorAxisLength',...
'MinorAxisLength', 'Eccentricity', 'EquivDiameter', ...
'Solidity', 'Extent', 8 );
```

```

% now convert it into array (matrix) while collecting
intensity data for each object:

%preallocate output array:
numobjects = size(ImStats, 1);
OData = zeros(numobjects, 14);
%now convert ImStats into array and add intensity data to
it
for k=1:numobjects
    OData(k, 1) = k;
    OData(k, 2) = ImStats(k).Centroid(1);
    OData(k, 3) = ImStats(k).Centroid(2);
    OData(k, 4) = ImStats(k).Area;
    OData(k, 5) = (ImStats(k).MajorAxisLength) /
(ImStats(k).MinorAxisLength);
    OData(k, 6) = ImStats(k).Eccentricity ;
    OData(k, 7) = ImStats(k).EquivDiameter;
    OData(k, 8) = ImStats(k).Solidity;
    OData(k, 9) = ImStats(k).Extent;

    % now collect and assign intensity parameters from
image I

    object_pixels = find( Ilabel == k);
    object_area = size(object_pixels, 1); %same as total
number of pixels in the object
    object_intensities = double(I(object_pixels)); % need
to convert to double to do math
    sorted_intensities = sort(object_intensities); % will
need to get median, 20% and 80% pixels
    total_intensity = sum(object_intensities, 1);
    avg_intensity = total_intensity / object_area;

```

```

        median_intensity = sorted_intensities( floor(
object_area/2 ) + 1 );
        pix20 = sorted_intensities( floor(object_area*0.2)+1 )
; %brightest pixel among dimmest 20%
        pix80 = sorted_intensities( floor(object_area*0.8)+1 )
;

        OData(k, 10) = total_intensity;
        OData(k, 11) = avg_intensity;
        OData(k, 12) = median_intensity;
        OData(k, 13) = pix20; %brightest pixel among dimmest
20%
        OData(k, 14) = pix80; %dimmest pixel among brightest
20%
    end %for

%===== end function
GetObjectsData()=====

function Imask = MaskDNA1(I);
% MaskDNA1 - generates binary mask for cell nuclei through
edge detection
% EV 1-22-99; 2-6-99; 2-10-99
% Imask = MaskDNA1(I)
% PARAMETERS
%   I - intensity image (grayscale)
% OUTPUT
%   Imask - BW image with objects from I
%
% For more details see Notebook Matlab_DNA_masking1_1-22-
99.doc
% Uses SDC Morphology Toolbox V0.7

```

```

if (nargin ~= 1)
    error('Wrong number of input parameters');
end
if (nargout ~= 1)
    error('Wrong number of output parameters: one output
argument should be provided');
end

Imask = edge(I, 'canny');
Imask = mm dil(Imask, mmsecross(1));
Imask = mmero ( mmc lohole(Imask,mmsecross(1)));
Imask = mm edgeoff(Imask, mmsecross(1));
% note that mm edgeoff this command removed FILLED OBJECTS
but not touching OUTLINES.
% these outlines can be removed by filtering:
Imask = medfilt2(Imask, [5 5]);

%=====end MaskDNA1 =====

```

Given the list of image files or montages of images as an input, this program creates an individual file for each image that contains the following quantitative measurements for all objects identified in the image:

- |   |                                    |
|---|------------------------------------|
| 1 - Index ("number" of an object);      | 8 - Solidity;                      |
| 2 - X coordinate of the center of mass; | 9 - Extent;                        |
| 3 - Y coordinate        "-";            | 10 - Total Intensity;              |
| 4 - Total Area (in pixels);             | 11 - Avg. Intensity;               |
| 5 - Ratio of MajorAxis/MinorAxis;       | 12 - Median Intensity;             |
| 6 - Eccentricity;                       | 13 - Intensity of 20% bright pixel |
| 7 - EquivDiameter;                      | 14 - Intensity of 80% bright pixel |

A fragment of an output for a single file, containing 9 images of cells stained for DNA and acquired with a 10x objective. A montage image that was used as a source to generate data in A is presented in Fig. 17.

The same program also summarizes measurements across many files and performs statistical analysis of the summary data. It creates a summary file with the following data:

- |   |                                    |
|---|------------------------------------|
| 1 - Image file number;                              |                                    |
| 2 - Average object Area (in pixels);                | 3 - STD (standard deviation) of 2; |
| 4 - Avg. of Ratio of MajorAxis/MinorAxis;           | 5 - STD of 4;                      |
| 6 - Avg. Eccentricity;                              | 7 - STD of 6;                      |
| 8 - Avg. EquivDiameter;                             | 9 - STD of 8;                      |
| 10 - Avg. of Solidity;                              | 11 - STD of 10;                    |
| 12 - Avg. of Extent;                                | 13 - STD of 11                     |
| 14 - Avg. of objects Total Intensity;               | 15 - STD of 14                     |
| 16 - Avg. of objects Avg Intensity;                 | 16 - STD of 15                     |
| 18 - Avg. of objects Median intensity;              | 19 - STD of 18                     |
| 20 - Avg. of objects intensity of 20% bright pixel; | 21 - STD of 19                     |
| 22 - Avg. of objects intensity of 80% bright pixel; | 23 - STD of 21                     |

An example of summary output obtained by running AnalyseDNA against 10 montage files also is shown in Appendix B.

A seventh experiment was conducted in order to use sequence analysis algorithms to analyze features of cell images. In this experiment, HeLa cells were treated for 24 hours with several different compounds, and then fixed, and stained with a fluorescent DNA dye. One image of these cells was acquired for each of the treatments and morphometric parameters and features were measured:

Resulting measurements were arranged into a string of numbers and reduced to a pseudo- nucleic acid sequence using following rules: At any given position in the sequence a number was substituted by "t" (a code for thymidine) if its value is among highest 25% of the values at the corresponding position in the data set, "g" if it is between 50% and 25%, "c" if it is between 75% and 50%, and "a" if it belongs to lowest

25% of values. Thus one descriptor or sequence was generated per treatment as illustrated in Fig. 18.

Resulting sequences were clustered using an AlignX module commercial software package Vector NTI (<http://informaxinc.com>), which uses a Neighbor Joining algorithm for sequence clustering.

The resulting dendrogram is presented in Fig 18. On the dendrogram the closest "leafs" correspond to the closest pseudo-sequences. Interestingly, compounds with similar mechanisms of action cluster together on the dendrogram. Another example of the generation of pseudo-sequences and clustering is shown in Fig. 19.

In some embodiments, techniques according to the present invention can provide tools for the later stages of drug development such as clinical trial design and patient management. The properties of known drugs such as clinical trial and patient response information will be used in a similar fashion as the pre-clinical information to provide predictions about the properties of novel compounds. Because the human cell is the locus of drug action, a database containing drug-cell interactions can be able to provide predictive information for this aspect of drug development.

Although the above has generally described the present invention according to specific systems, the present invention has a much broader range of applicability. In particular, the present invention is not limited to a particular kind of data about a cell, but can be applied to virtually any cellular data where an understanding about the workings of the cell is desired. Thus, in some embodiments, the techniques of the present invention could provide information about many different types or groups of cells, substances, and genetic processes of all kinds. Of course, one of ordinary skill in the art would recognize other variations, modifications, and alternatives.



## CLAIMS

*what is claimed is:*

1. A method of characterizing a cell using one or more markers associated with components of the cell, the method comprising:
  - receiving images for each of the one or more markers;
  - for at least one of the images, determining one or more descriptors using the associated marker, wherein each descriptor characterizes the marker in a particular morphological or statistical manner; and
  - classifying the cell into three or more biologically relevant classes using at least one of the markers and associated descriptors.
2. The method of claim 1, wherein the cell has been exposed to a particular stimulus prior to imaging.
3. The method of claim 2, wherein the stimulus includes at least of exposure to a chemical agent, exposure to a biological agent, exposure to radiation, and combinations thereof.
4. The method of claim 1, wherein at least one of the components of the cell is selected from the group consisting of DNA, Golgi, cytoskeletal proteins, and combinations thereof.
5. The method of claim 1, wherein one of the markers labels DNA.
6. The method of claim 1, wherein one of the markers labels a Golgi component.
7. The method of claim 1, wherein one of the markers labels a cytoskeletal protein.
8. The method of claim 7, wherein the cytoskeletal protein is tubulin.
9. The method of claim 1, wherein the markers include at least DNA, a Golgi component and tubulin.
10. The method of claim 1, wherein prior to receiving images, the cell was treated with a first marker that binds to a first cell component and emits a signal in proportion to the concentration of the first marker.

11. The method of claim 10, wherein prior to receiving the images, the cell was treated with a second marker that binds to a second cell component and emits a signal in proportion to the concentration of the second marker.
12. The method of claim 1, wherein at least one of the one or more descriptors characterizes an intensity histogram of the associated image.
13. The method of claim 1, wherein at least one of the one of more descriptors characterizes a morphologic property of a cellular component within the cell.
14. The method of claim 1, wherein classifying the cell uses at least DNA and one other marker.
15. The method of claim 14, wherein classifying the cell uses at least DNA and a Golgi component.
16. The method of claim 14, wherein classifying the cell uses at least DNA and a cytoskeletal protein.
17. A computer program product comprising a machine readable medium on which is provided program instructions for characterizing a cell using one or more markers associated with components of the cell, the program instructions comprising:
  - program code for receiving images for each of the one or more markers;
  - program code for determining one or more descriptors, for at least one of the images, using the associated marker, wherein each descriptor characterizes the marker in a particular morphological or statistical manner; and
  - program code for classifying the cell into three or more biologically relevant classes using at least one of the markers and associated descriptors.
18. The computer program product of claim 17, wherein the cell has been exposed to a particular stimulus prior to imaging.
19. The computer program product of claim 18, wherein the stimulus includes at least of exposure to a chemical agent, exposure to a biological agent, exposure to radiation, and combinations thereof.

20. The computer program product of claim 17, wherein at least one of the components of the cell is selected from the group consisting of DNA, Golgi, cytoskeletal proteins, and combinations thereof.
21. The computer program product of claim 17, wherein one of the markers labels DNA.
22. The computer program product of claim 17, wherein one of the markers labels a Golgi component.
23. The computer program product of claim 17, wherein one of the markers labels a cytoskeletal protein.
24. The computer program product of claim 23, wherein the cytoskeletal protein is tubulin.
25. The computer program product of claim 17, wherein the markers include at least DNA, a Golgi component and tubulin.
26. The computer program product of claim 17, wherein the cell was treated with a first marker that binds to a first cell component and emits a signal in proportion to the concentration of the first marker.
27. The computer program product of claim 26, wherein the cell was treated with a second marker that binds to a second cell component and emits a signal in proportion to the concentration of the second marker.
28. The computer program product of claim 17, wherein the program code for determining one or more descriptors comprises program code for determining at least one descriptor that characterizes an intensity histogram of the associated image.
29. The computer program product of claim 17, wherein the program code for determining one or more descriptors comprises program code for determining at least one descriptor that characterizes a morphologic property of a cellular component within the cell.
30. The computer program product of claim 17, wherein the program code for classifying the cell uses at least DNA and one other marker.

31. The computer program product of claim 30, wherein the program code for classifying the cell uses at least DNA and a Golgi component.
32. The computer program product of claim 30, wherein the program code for classifying the cell uses at least DNA and a cytoskeletal protein.
33. A method of characterizing cellular activity of a compound, the method comprising:  
receiving one or more images of cells having a known cellular activity as a result of one or more of the following: a genetic or epigenetic modification, treatment with a selected compound known to impart the cellular activity, and treatment with a plurality of compounds imparting varying levels of the known cellular activity;  
quantitatively characterizing phenotypic attributes of the one or more images of the cells to thereby produce a target phenotype for the known cellular activity;  
receiving one or more images of other cells treated with a compound to be characterized;  
quantitatively characterizing phenotypic attributes of the one or more images of the other cells to thereby produce a second phenotype associated with the compound to be characterized; and  
comparing the target phenotype and the second phenotype to determine whether the compound to be characterized possesses the known cellular activity.
34. The method of claim 33, wherein the known cellular activity results from treatment with a selected compound.
35. The method of claim 33, wherein the known cellular activity results from a genetic or epigenetic treatment.
36. The method of claim 33, wherein the known cellular activity results from treatment with a plurality of compounds imparting varying levels of the known cellular activity.
37. The method of claim 33, wherein the known cellular activity is at least of mechanism of action and toxicity.
38. The method of claim 37, wherein the mechanism of action is a mechanism of action for a cancer.
39. The method of claim 33, wherein the one or more images of cells having a known cellular activity are obtained from at least two different cell lines.
40. The method of claim 33, wherein the phenotypic attributes of the cells comprise attributes of one more cell markers.
41. The method of claim 40, wherein the markers are markers for one or more cellular organelles.

42. The method of claim 33, wherein the target phenotype comprises a fingerprint comprised of multiple scalar values of the phenotypic attributes.
43. The method of claim 33, further comprising comparing the target phenotype to second phenotypes for a plurality of other compounds.
44. The method of claim 43, further comprising ranking the plurality of other compounds based on degree of similarity to the target phenotype.
45. The method of claim 35, wherein the genetic or epigenetic modification comprises at least one genetically knocking out a gene, underexpressing the gene, overexpressing the gene, and expressing the gene in a non-native state.
46. The method of claim 36, wherein the plurality of compounds comprises positive and negative biochemical hits.
47. The method of claim 46, further comprising ranking the positive and negative hits based upon their interaction with the target.
48. A computer program product comprising a machine readable medium on which is provided program instructions for characterizing cellular activity of a compound, the program instructions comprising:  
    program code for receiving one or more images of cells having a **known cellular activity** as a result of one or more of the following: a genetic or epigenetic modification, treatment with a selected compound known to impart the cellular activity, and treatment with a plurality of compounds imparting varying levels of the known cellular activity;  
    program code for quantitatively characterizing phenotypic attributes of the one or more images of the cells to thereby produce a target phenotype for the known cellular activity;  
    program code for receiving one or more images of other cells treated with a **compound to be characterized**;  
    program code for quantitatively characterizing phenotypic attributes of the one or more images of the other cells to thereby produce a **second phenotype** associated with the compound to be characterized; and  
    program code for comparing the target phenotype and the second phenotype to determine whether the compound to be characterized possesses the known cellular activity.
49. The computer program product of claim 48, wherein the known cellular activity results from treatment with a selected compound.
50. The computer program product of claim 48, wherein the known cellular activity results from a genetic or epigenetic treatment.
51. The computer program product of claim 48, wherein the known cellular activity results from treatment with a plurality of compounds imparting varying levels of the known cellular activity.

52. The computer program product of claim 48, wherein the known cellular activity is at least of mechanism of action and toxicity.
53. The computer program product of claim 52, wherein the mechanism of action is a mechanism of action for a cancer.
54. The computer program product of claim 48, wherein the one or more images of cells having a known cellular activity are obtained from at least two different cell lines.
55. The computer program product of claim 48, wherein the phenotypic attributes of the cells comprise attributes of one more cell markers.
56. The computer program product of claim 55, wherein the markers are markers for one or more cellular organelles.
57. The computer program product of claim 48, wherein the target phenotype comprises a fingerprint comprised of multiple scalar values of the phenotypic attributes.
58. The computer program product of claim 48, further comprising program code for comparing the target phenotype to second phenotypes for a plurality of other compounds.
59. The computer program product of claim 58, further comprising program code for ranking the plurality of other compounds based on degree of similarity to the target phenotype.
60. The computer program product of claim 50, wherein the genetic or epigenetic modification comprises at least one genetically knocking out a gene, underexpressing the gene, overexpressing the gene, and expressing the gene in a non-native state.
61. The computer program product of claim 51, wherein the plurality of compounds comprises positive and negative biochemical hits.
62. The computer program product of claim 61, further comprising program code for ranking the positive and negative hits based upon their interaction with the target.

1/37

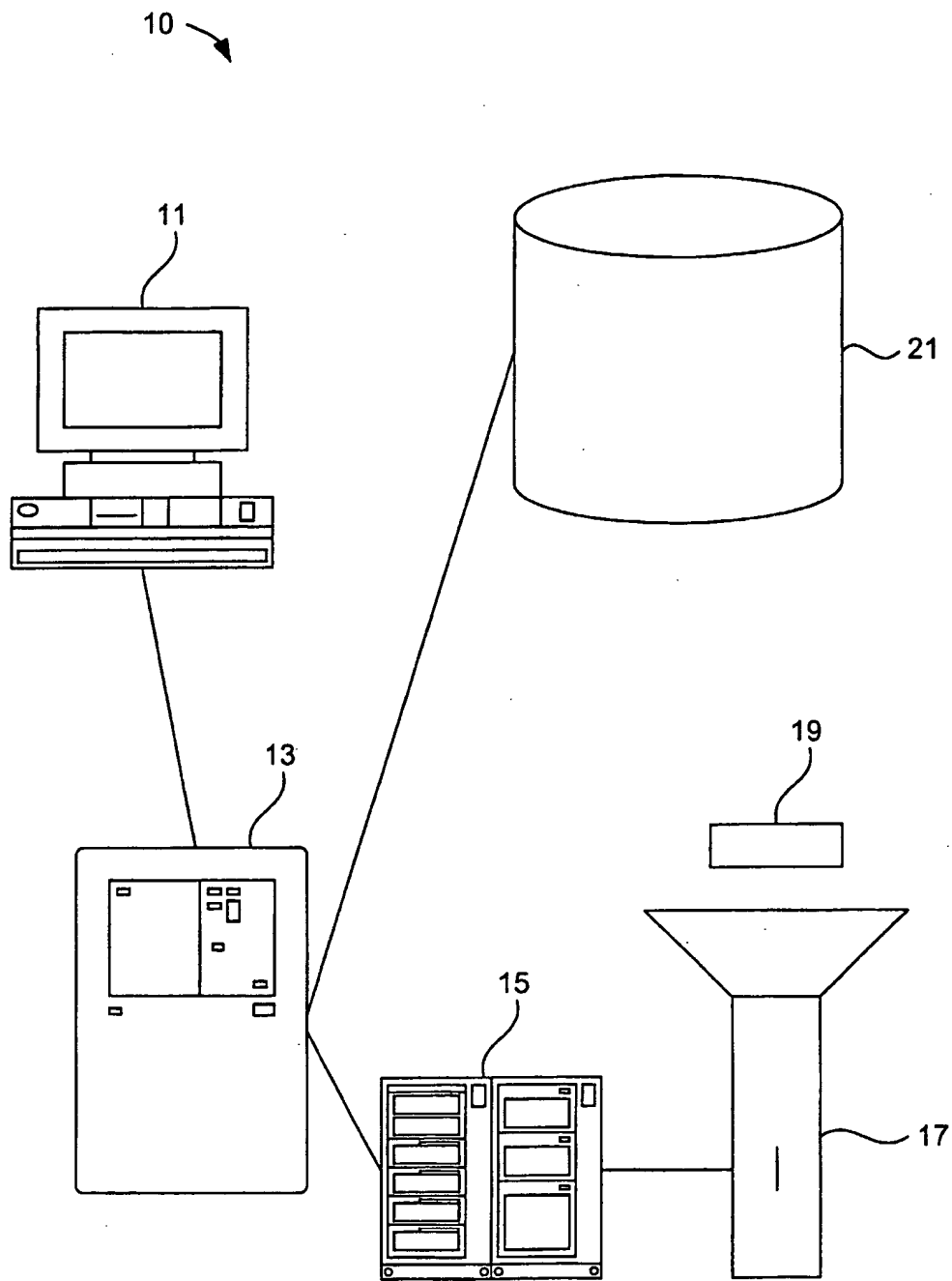


FIG. 1

2/37

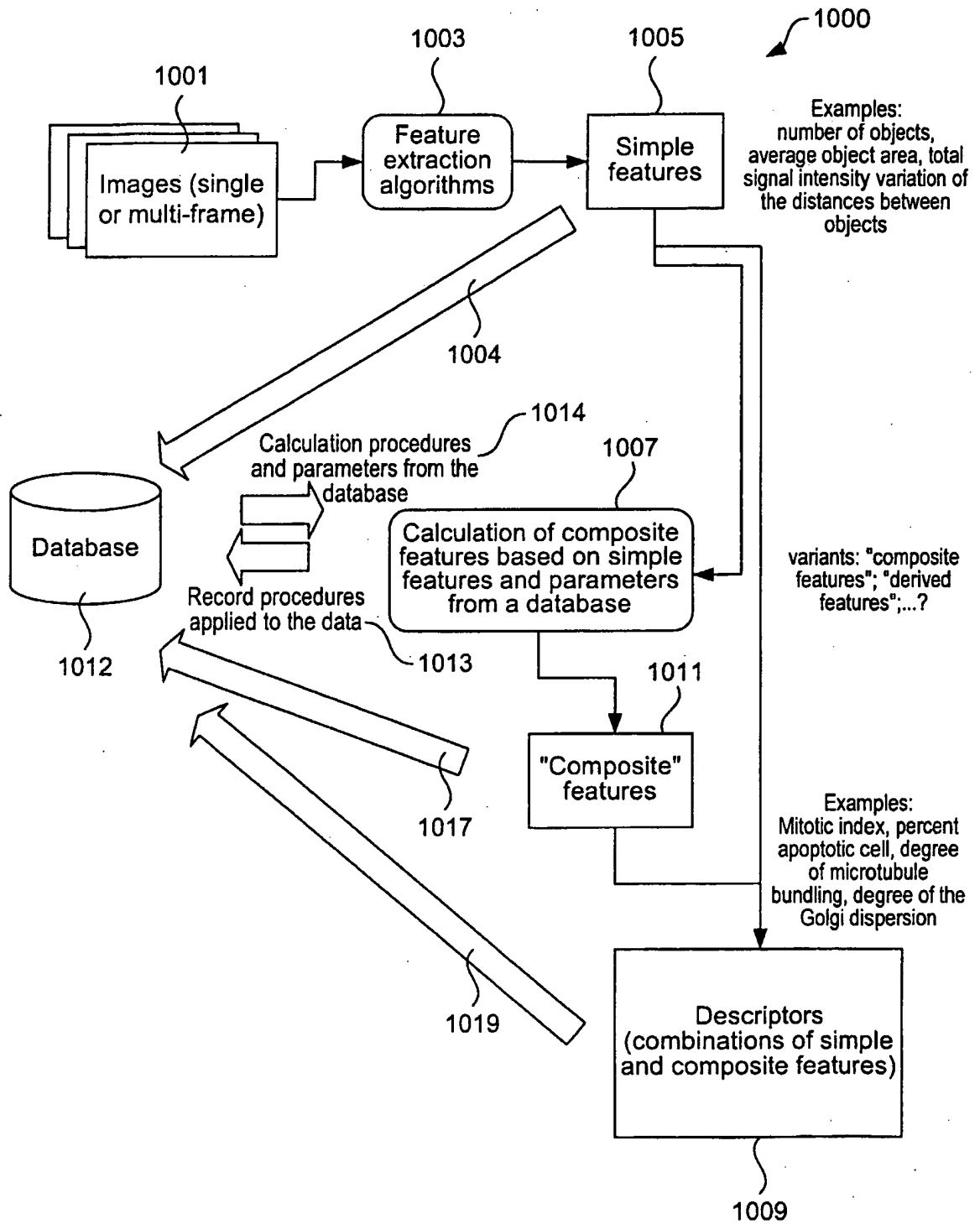


FIG. 1A



3/37

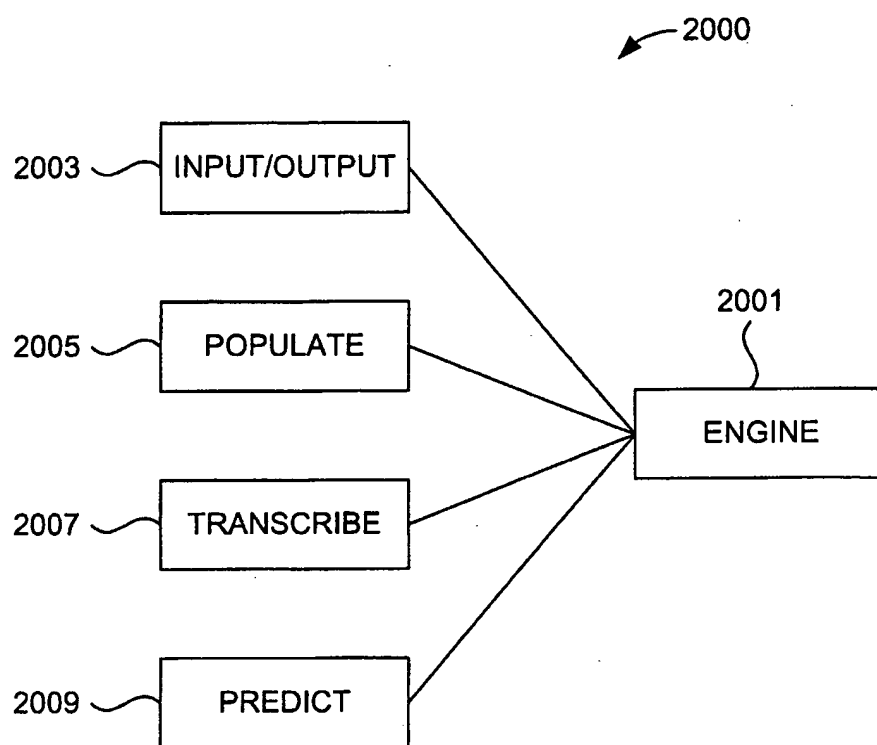


FIG. 1B

4/37

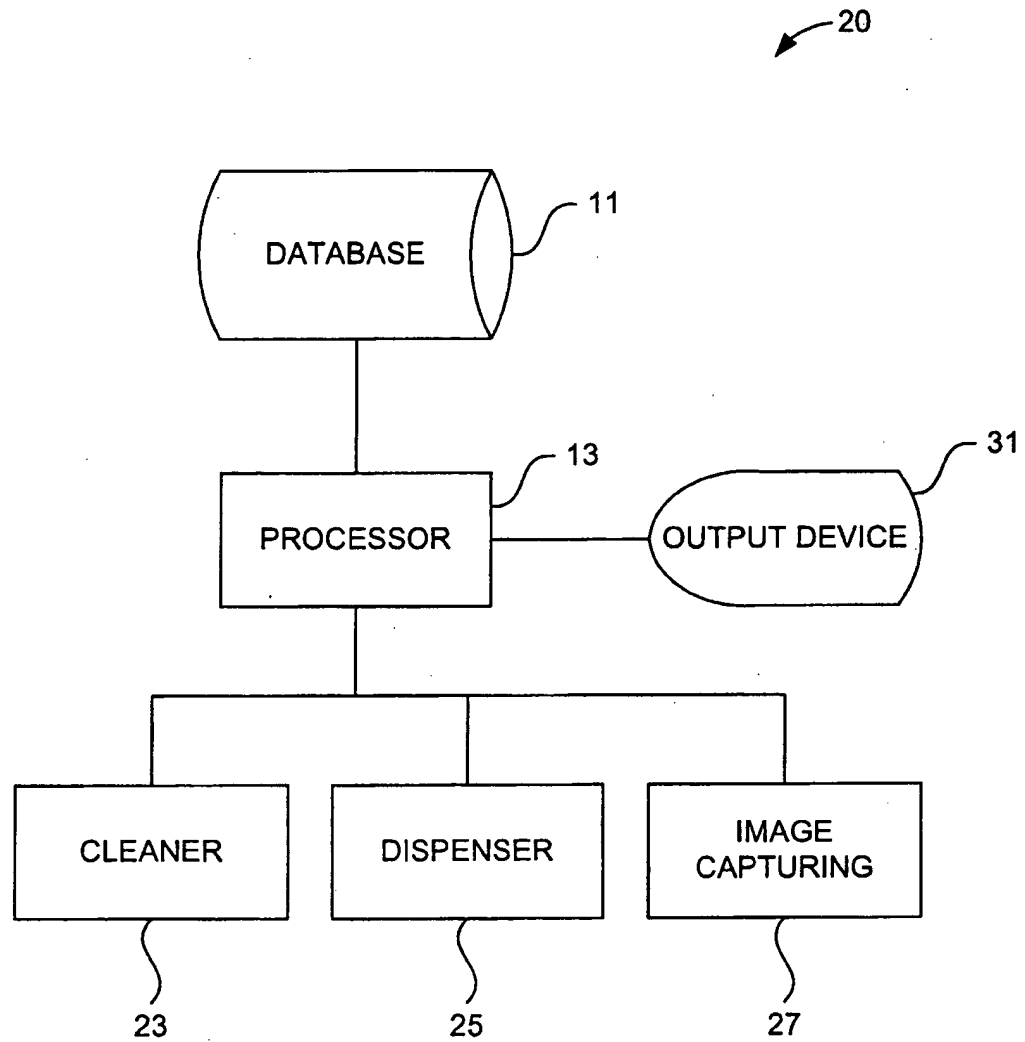


FIG. 2

5/37

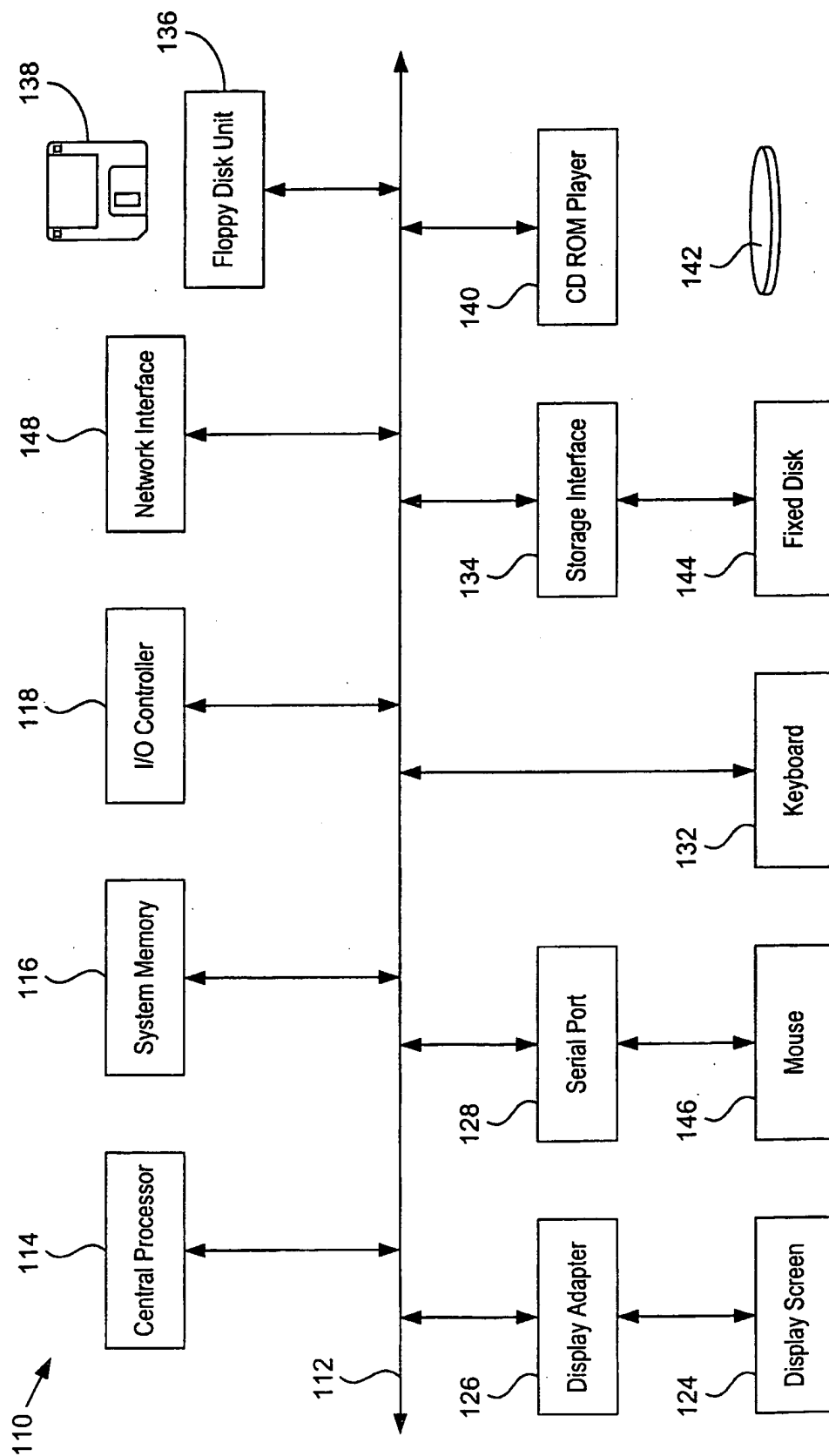


FIG. 3

6/37

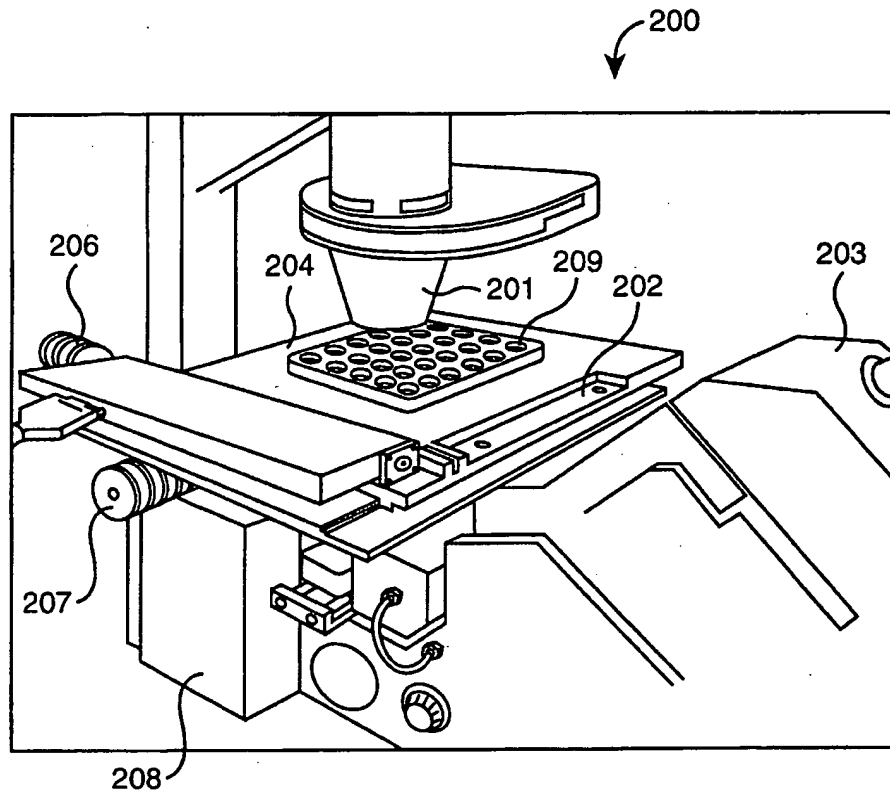


FIG. 4

7/37

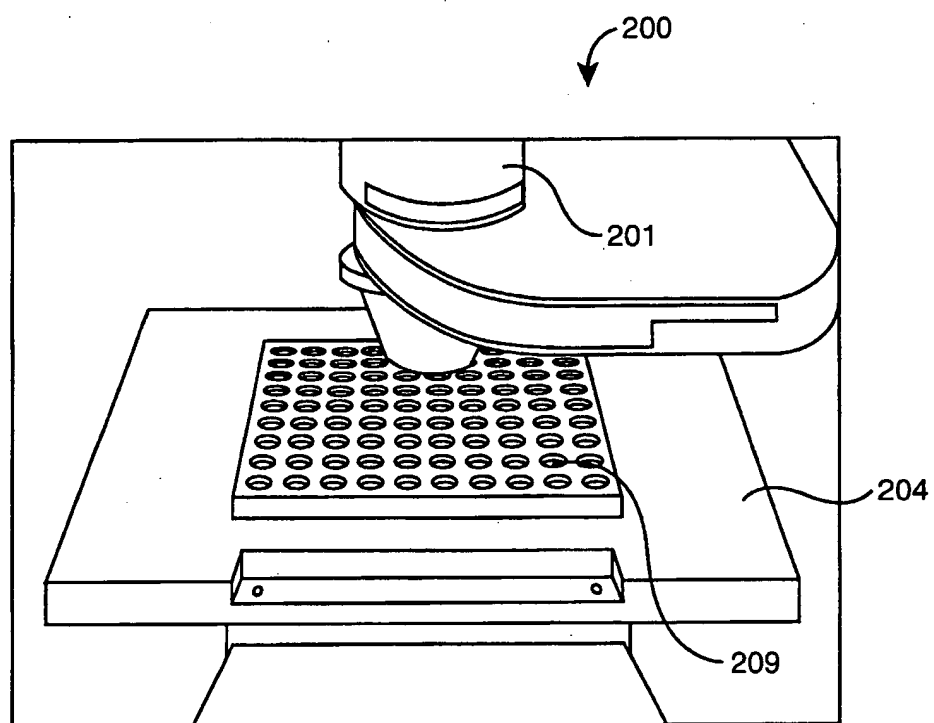


FIG. 5

8/37

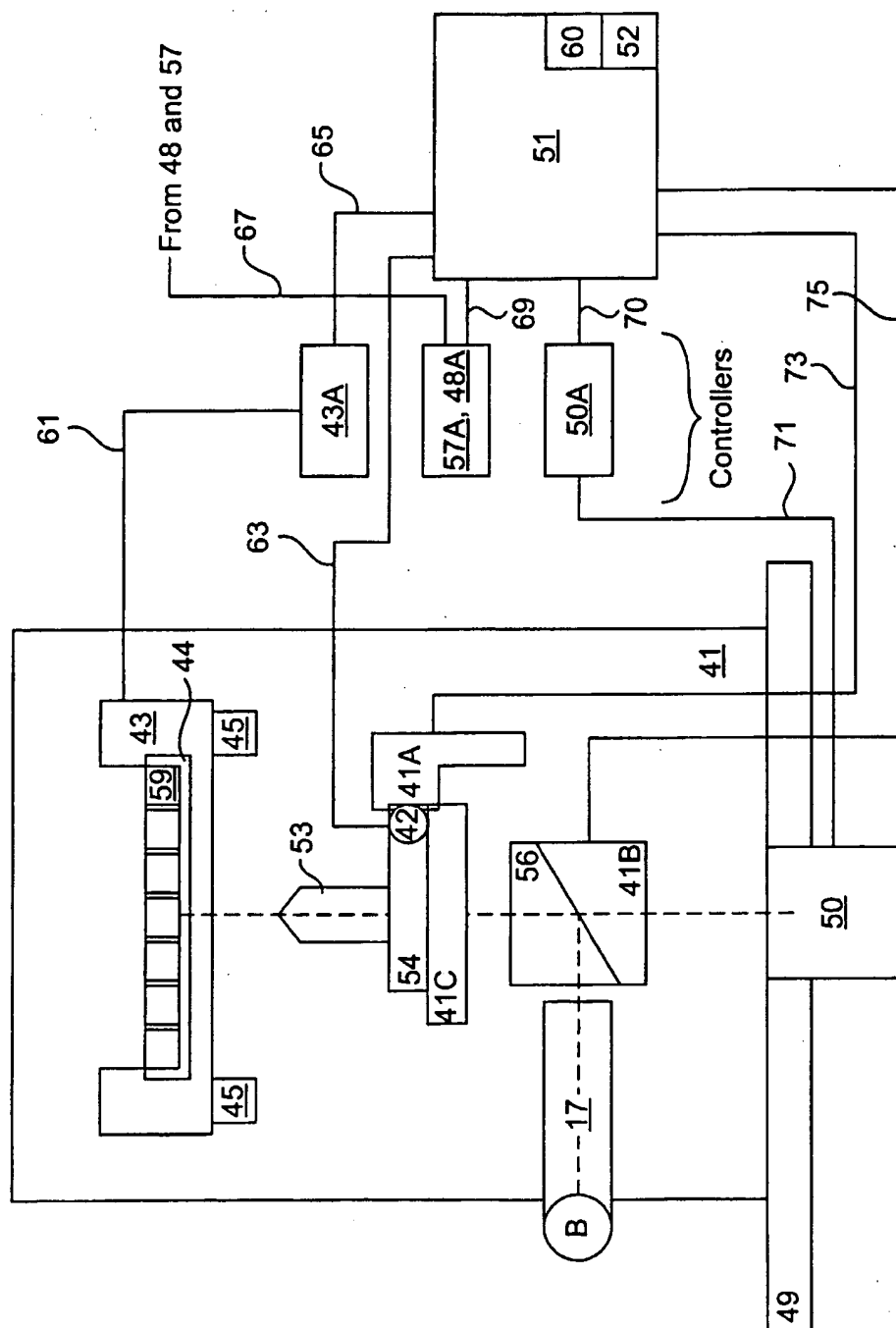


FIG. 5A

9/37

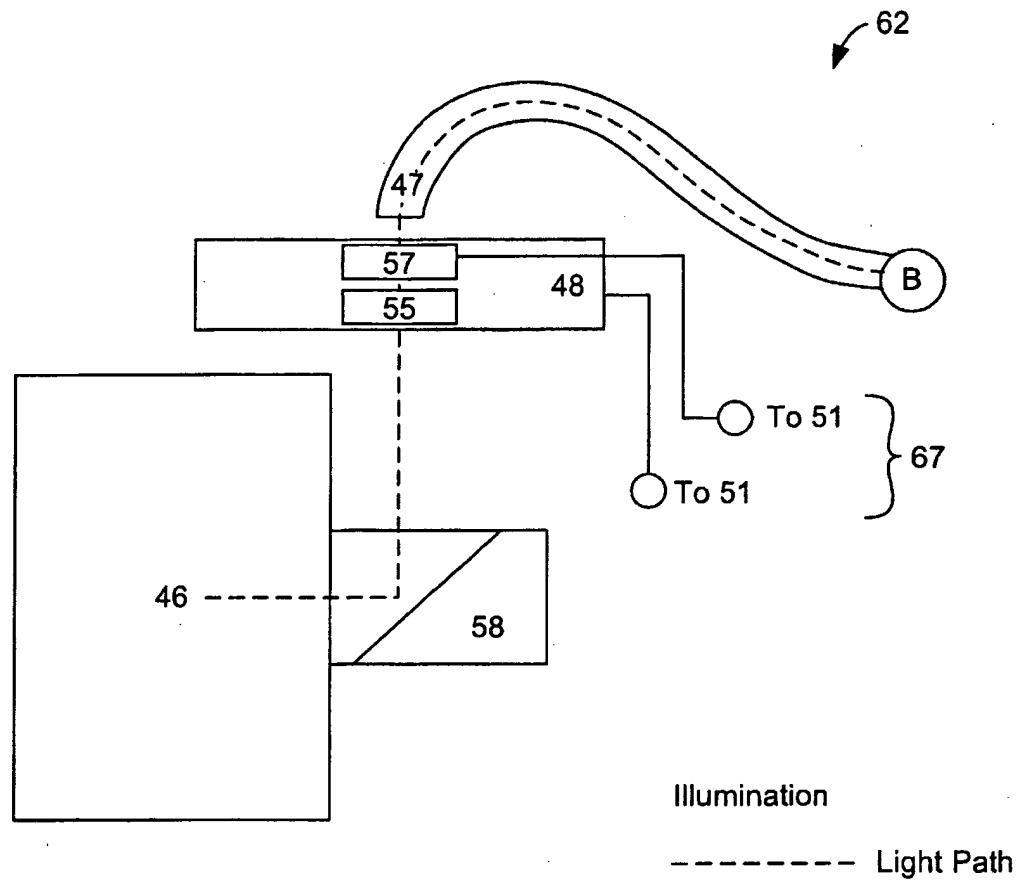


FIG. 5B

10/37

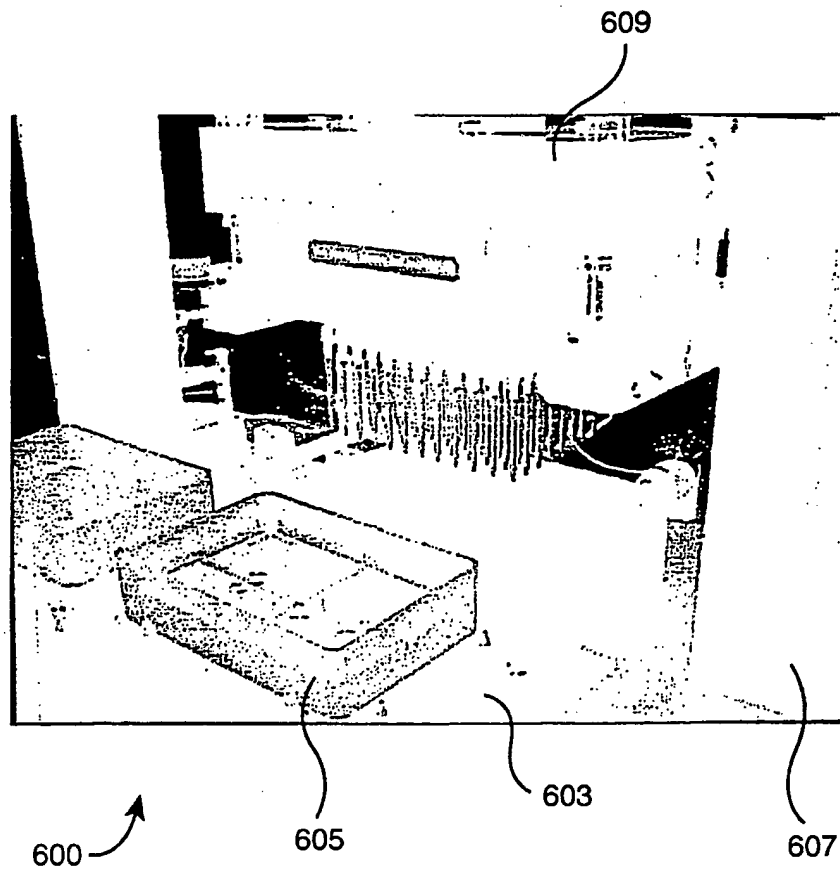


FIG. 6



11/37

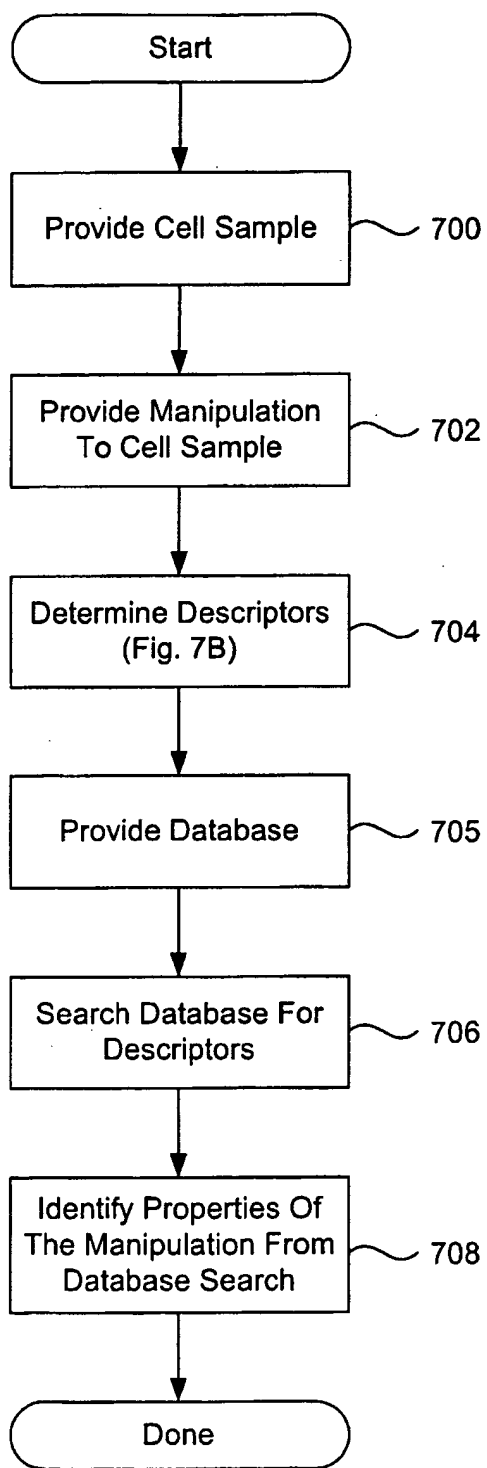
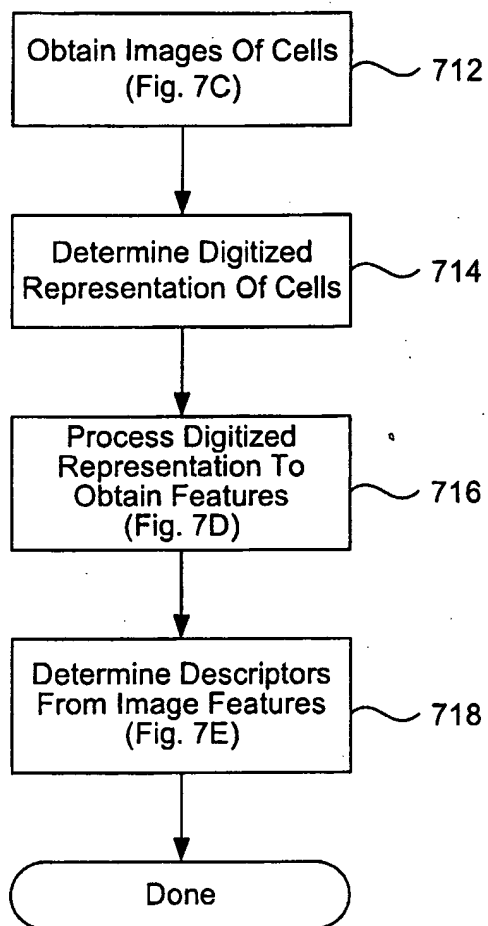


FIG. 7A

12/37



**FIG. 7B**  
Step 704 of Fig. 7A

13/37

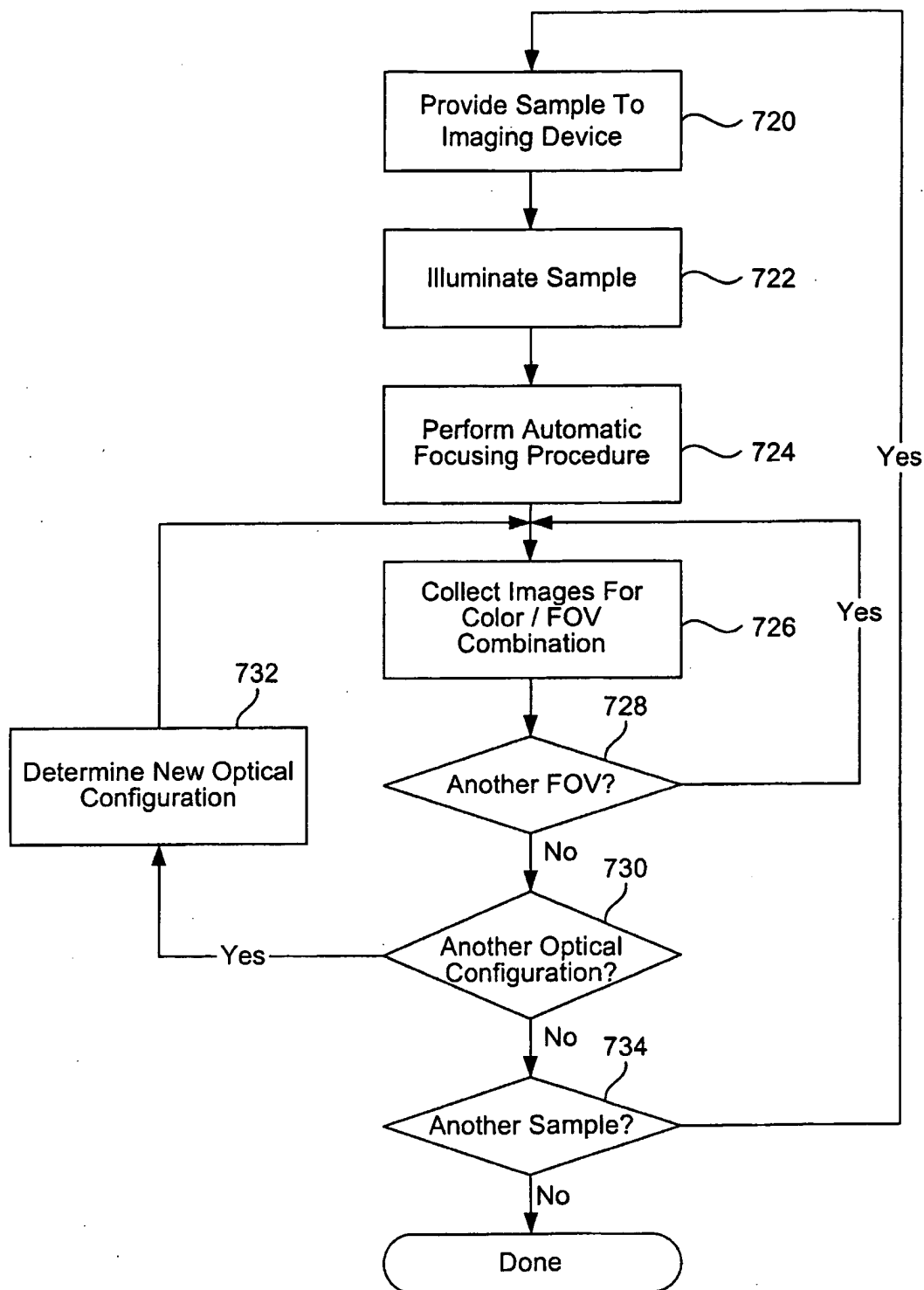
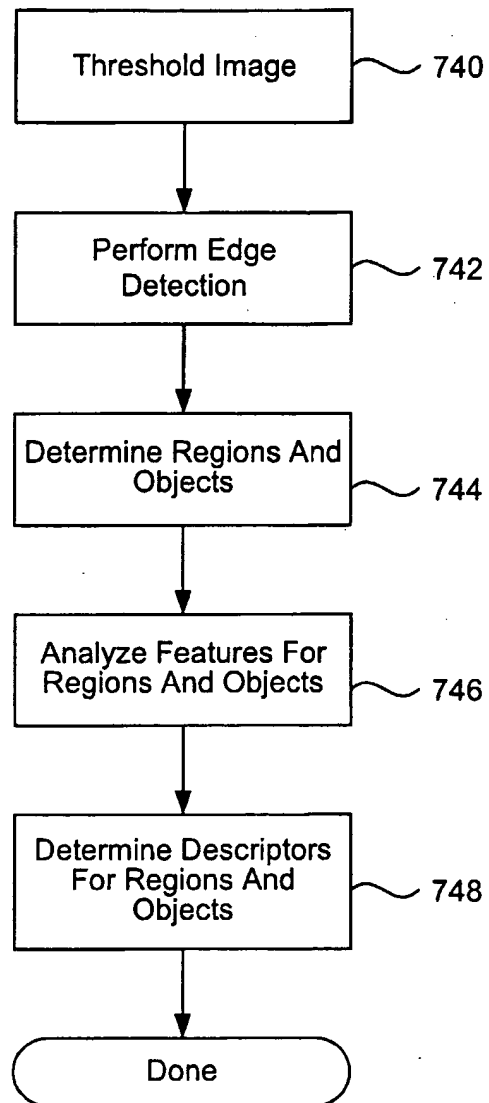


FIG. 7C  
Step 714 of Fig. 7B

14/37



**FIG. 7D**  
Step 716 of Fig. 7B

15/37

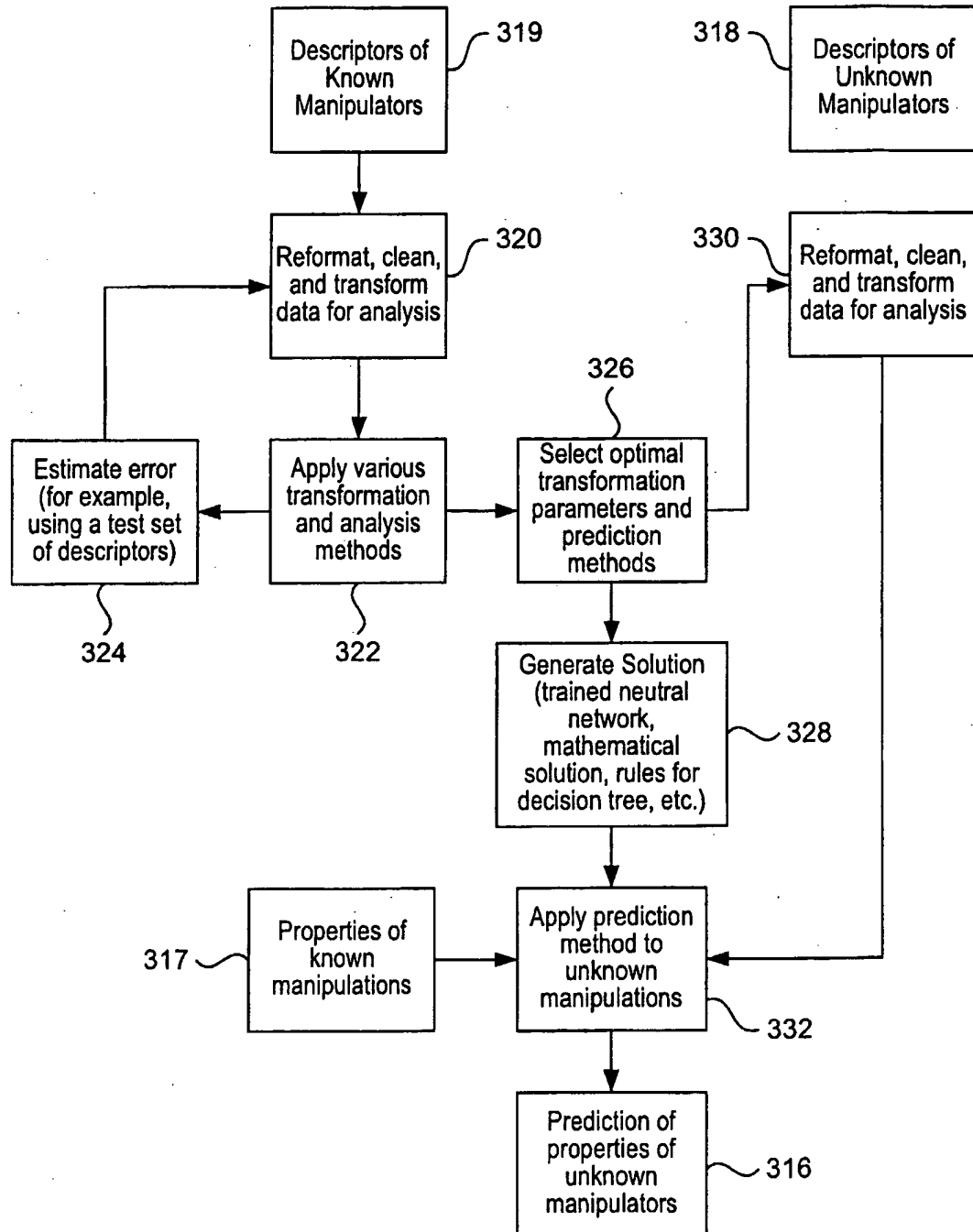


FIG. 7E

16/37

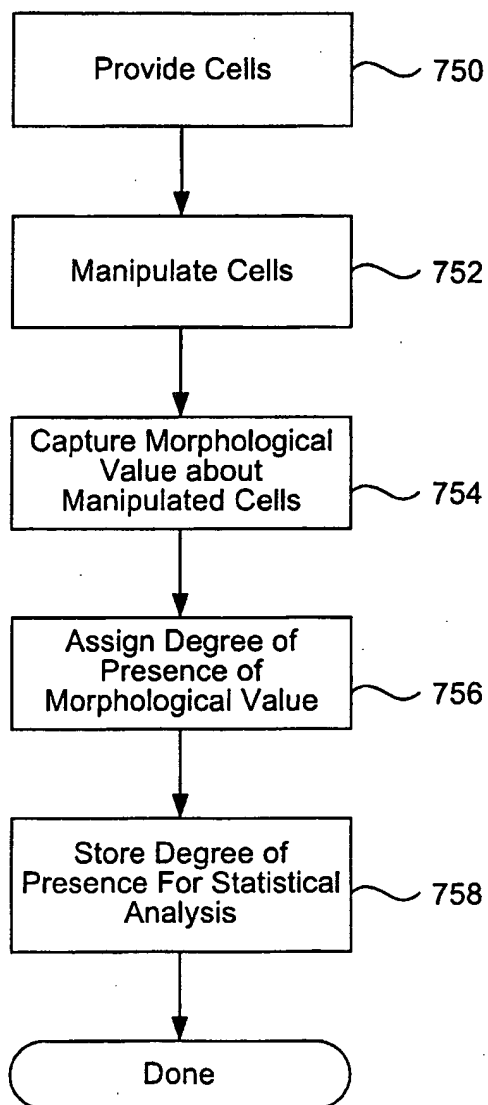


FIG. 7F

17/37

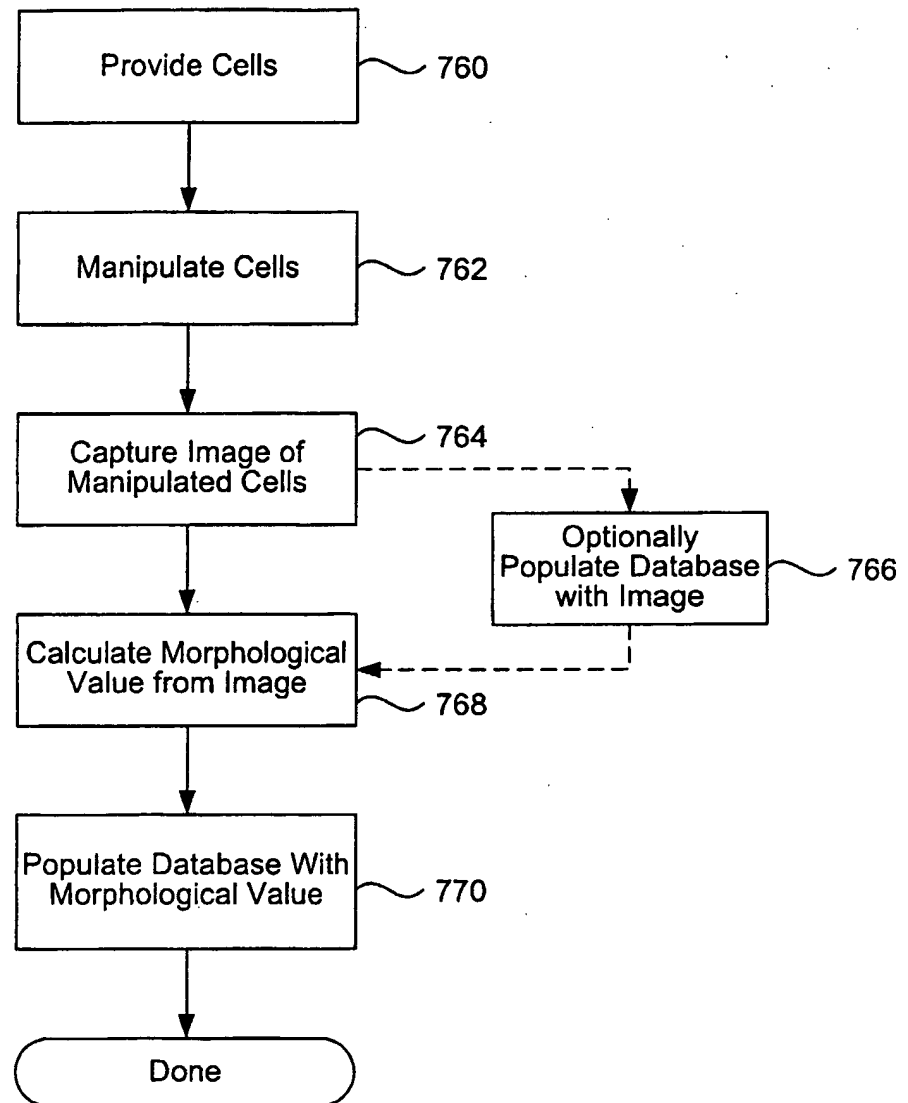


FIG. 7G

18/37

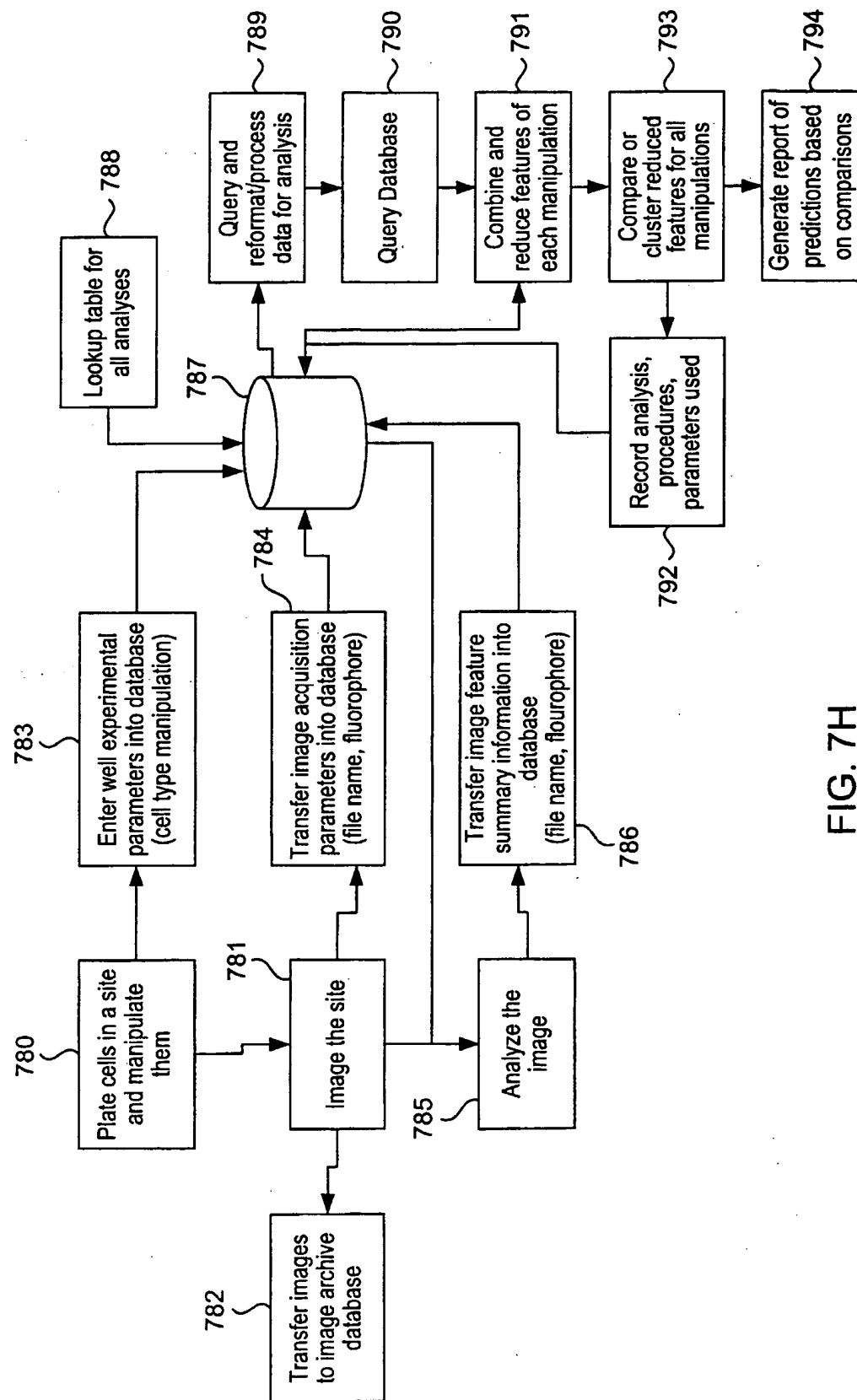


FIG. 7H



19/37

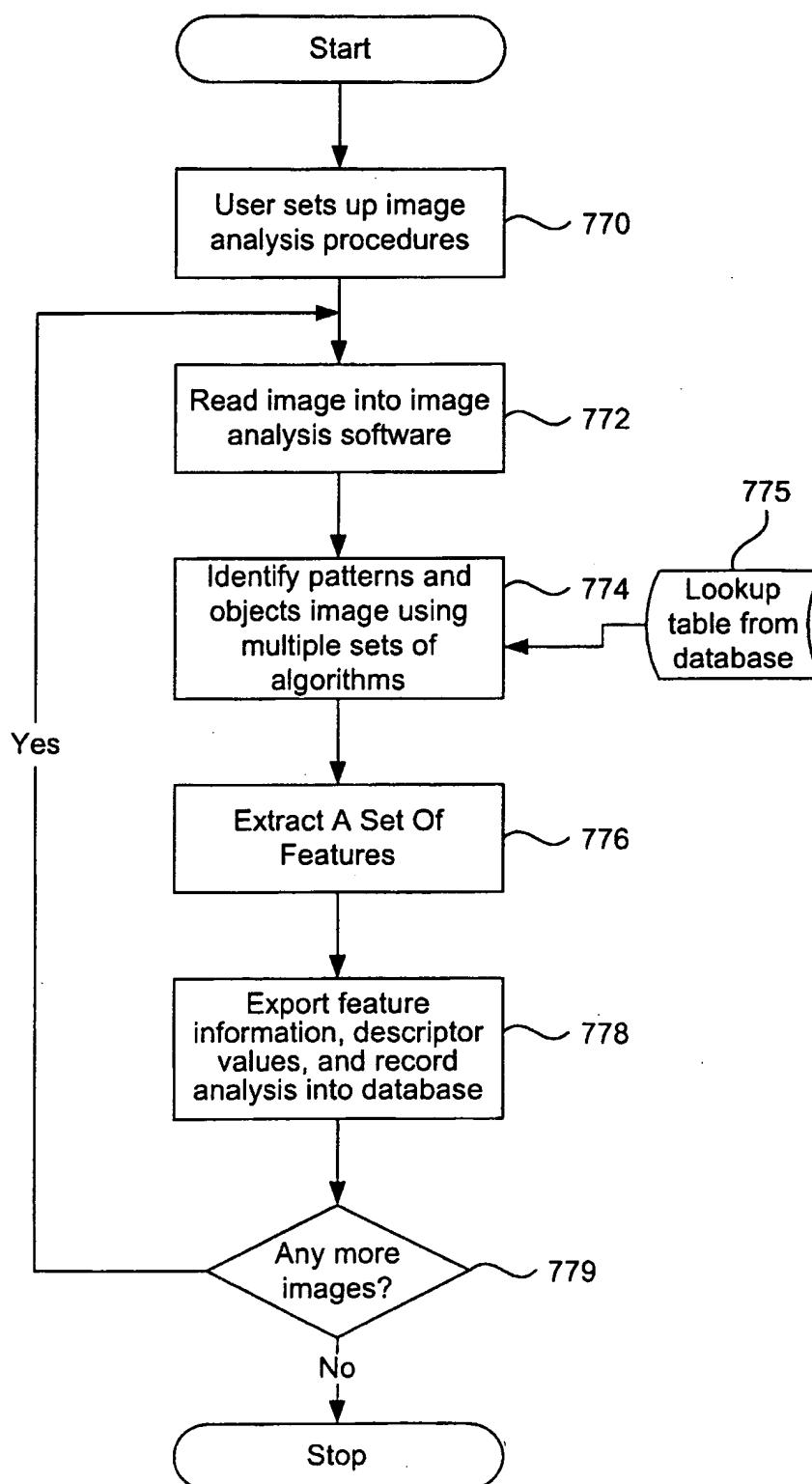


FIG. 71

20/37

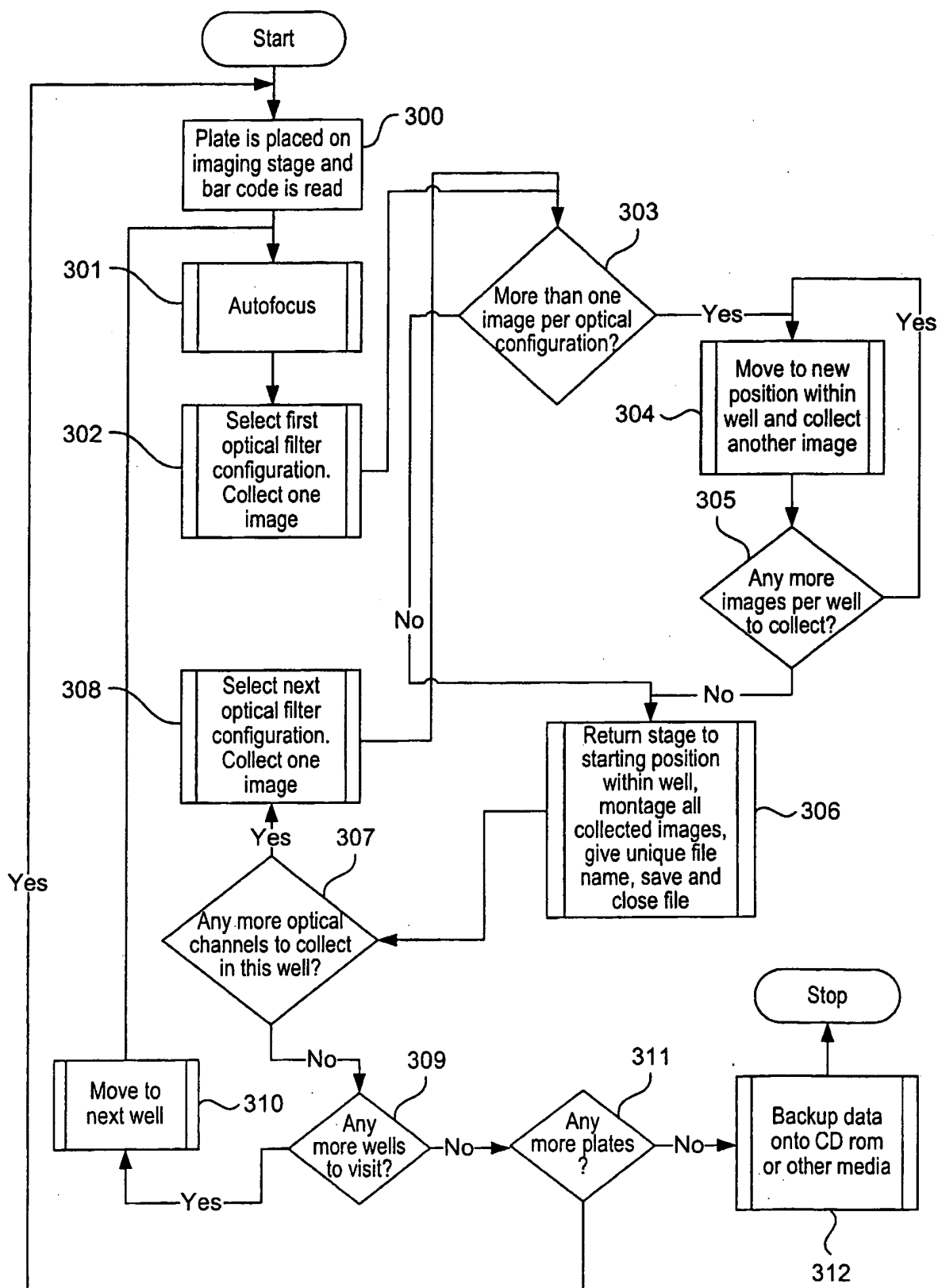


FIG. 7J

21/37

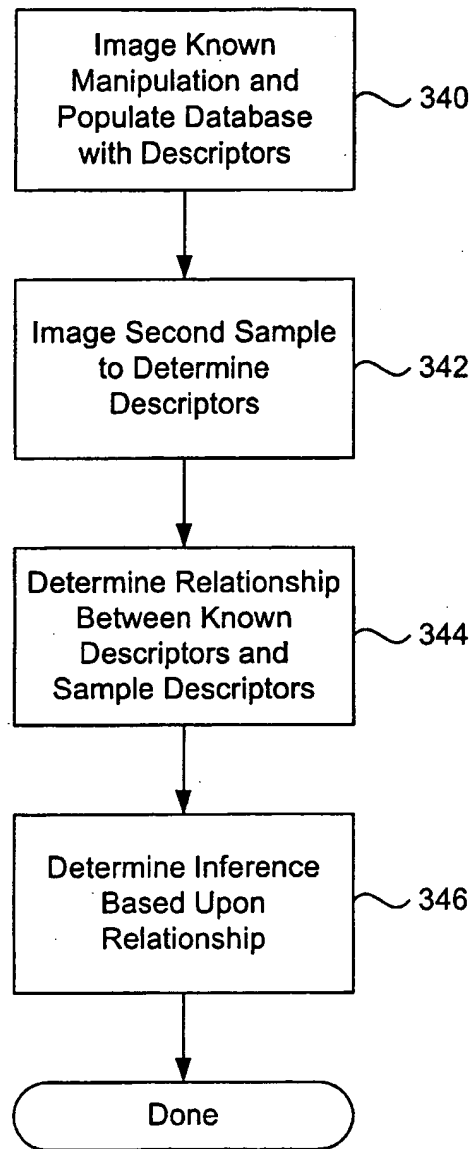


FIG. 7K

22/37

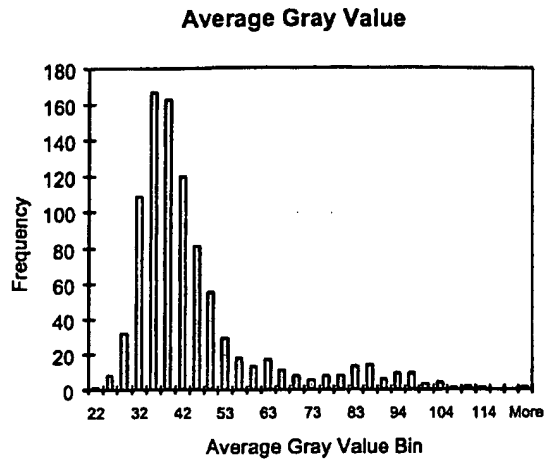


FIG. 8A

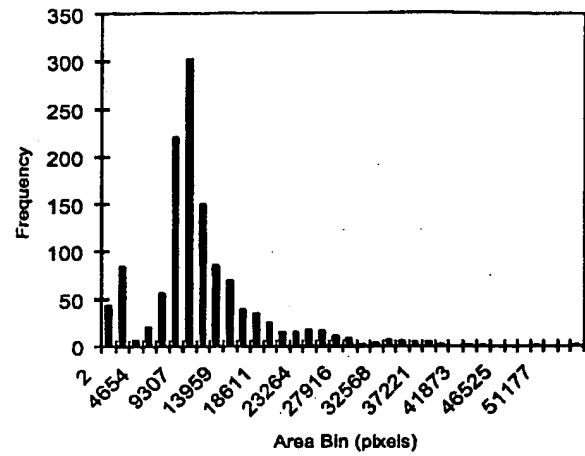


FIG. 8B

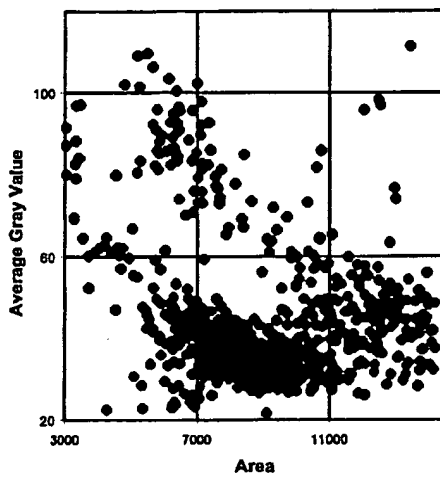


FIG. 8C

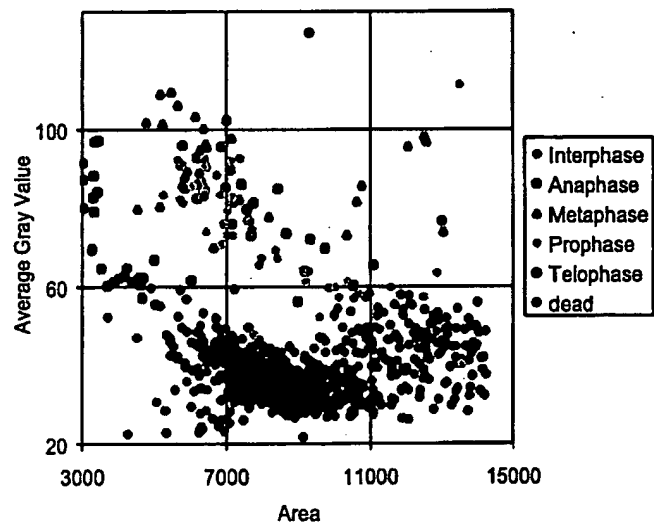


FIG. 8D

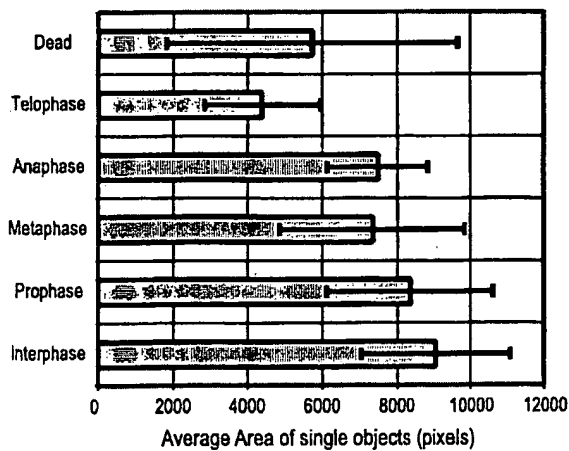


FIG. 8E

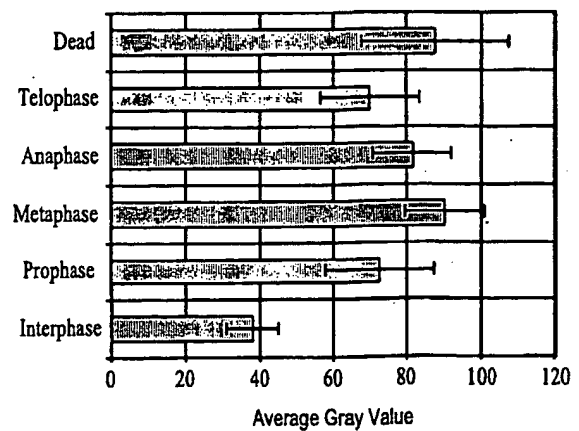


FIG. 8F

23/37

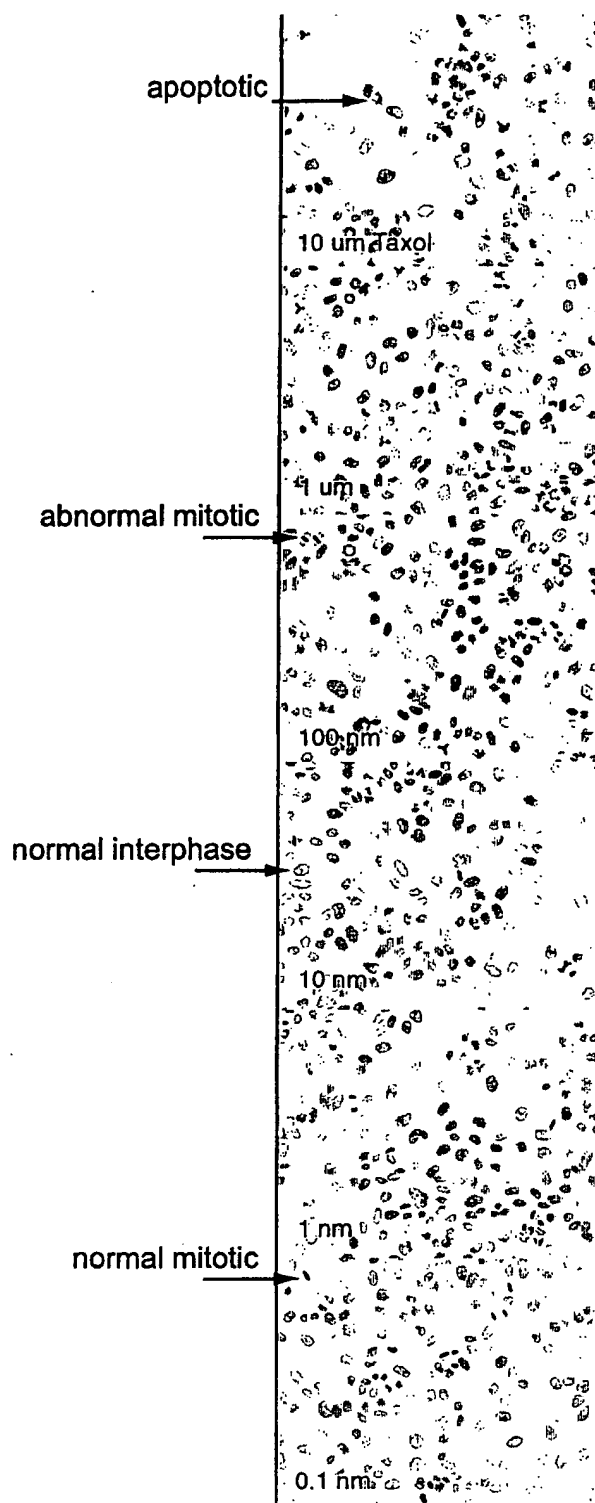


FIG. 9

24/37

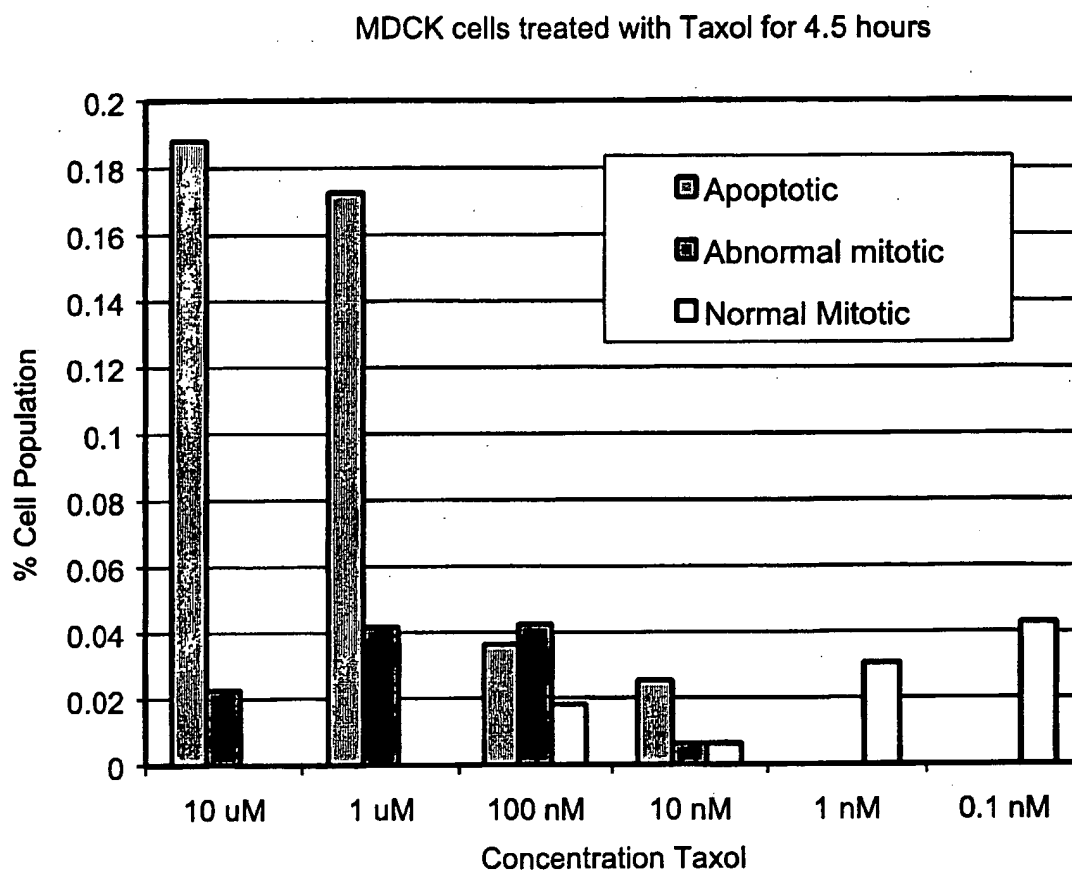


FIG. 10

25/37

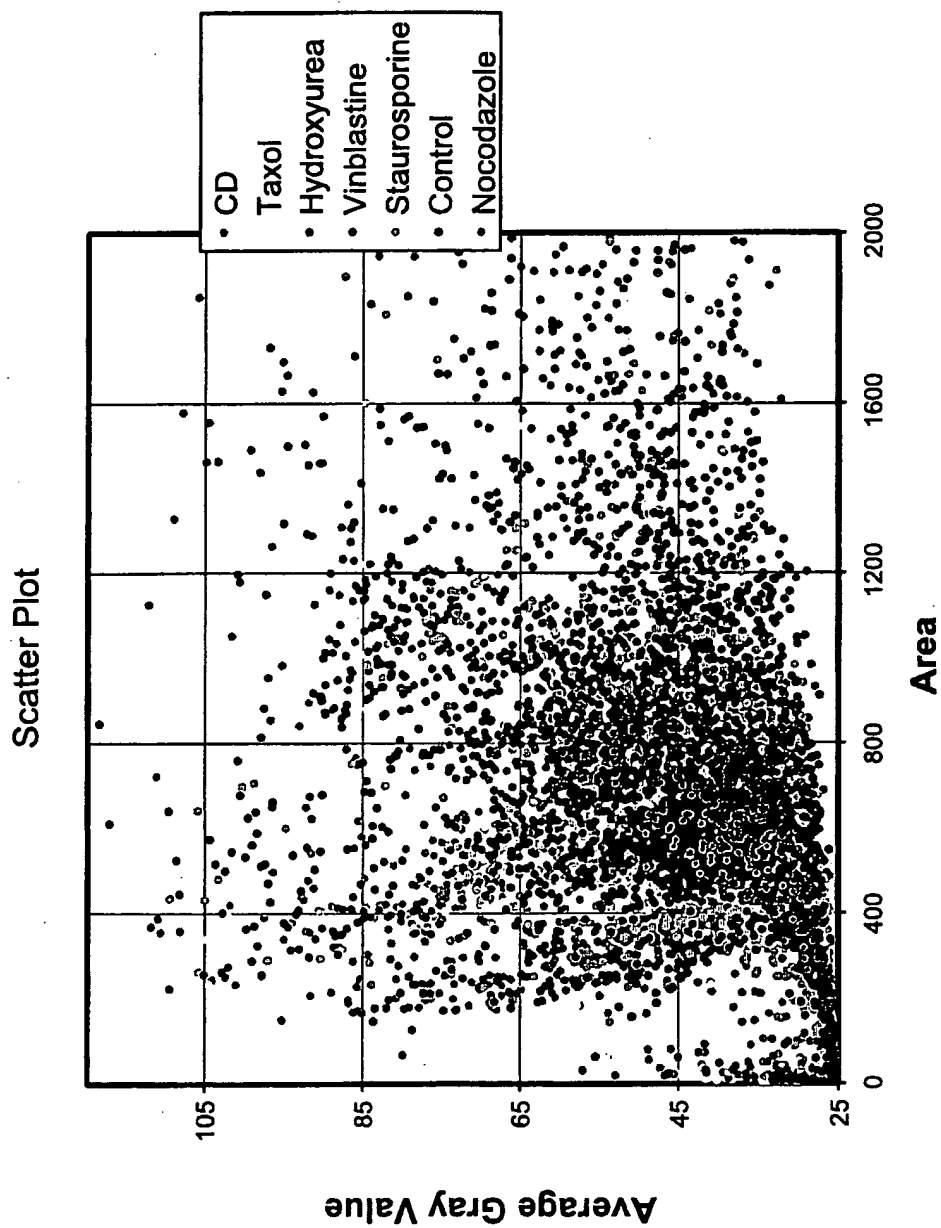


FIG. 11

26/37

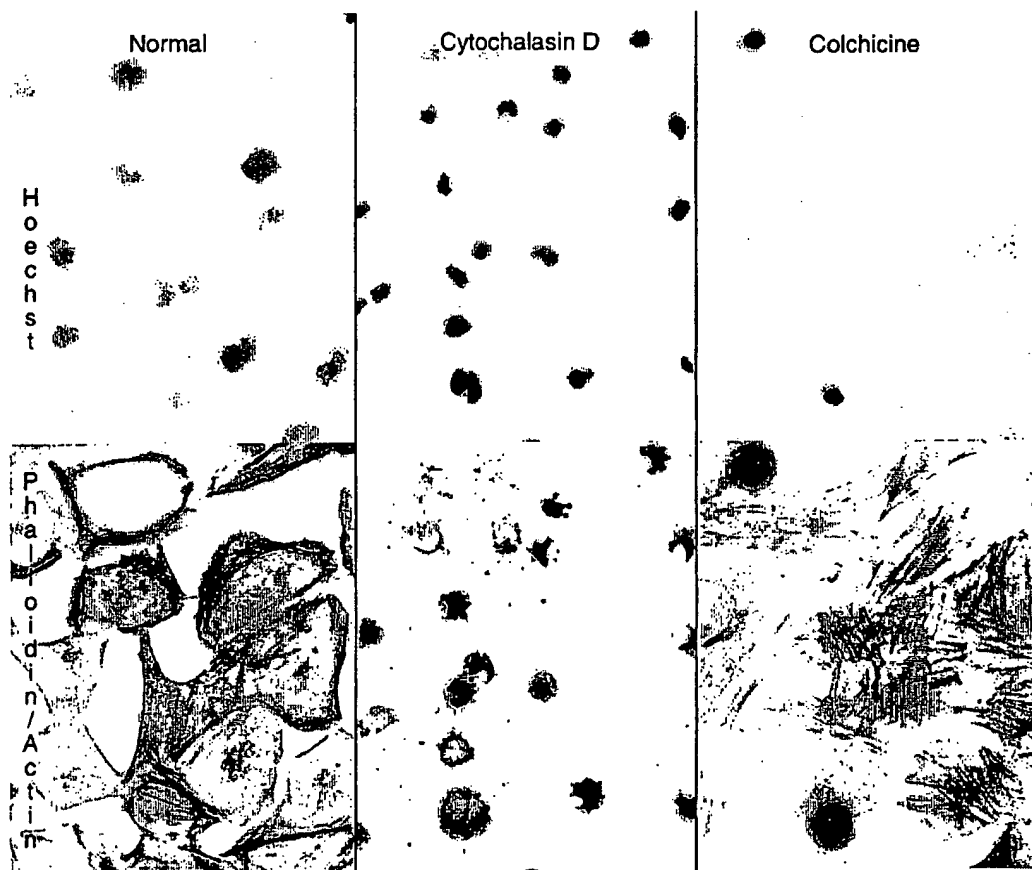


FIG. 12



27/37

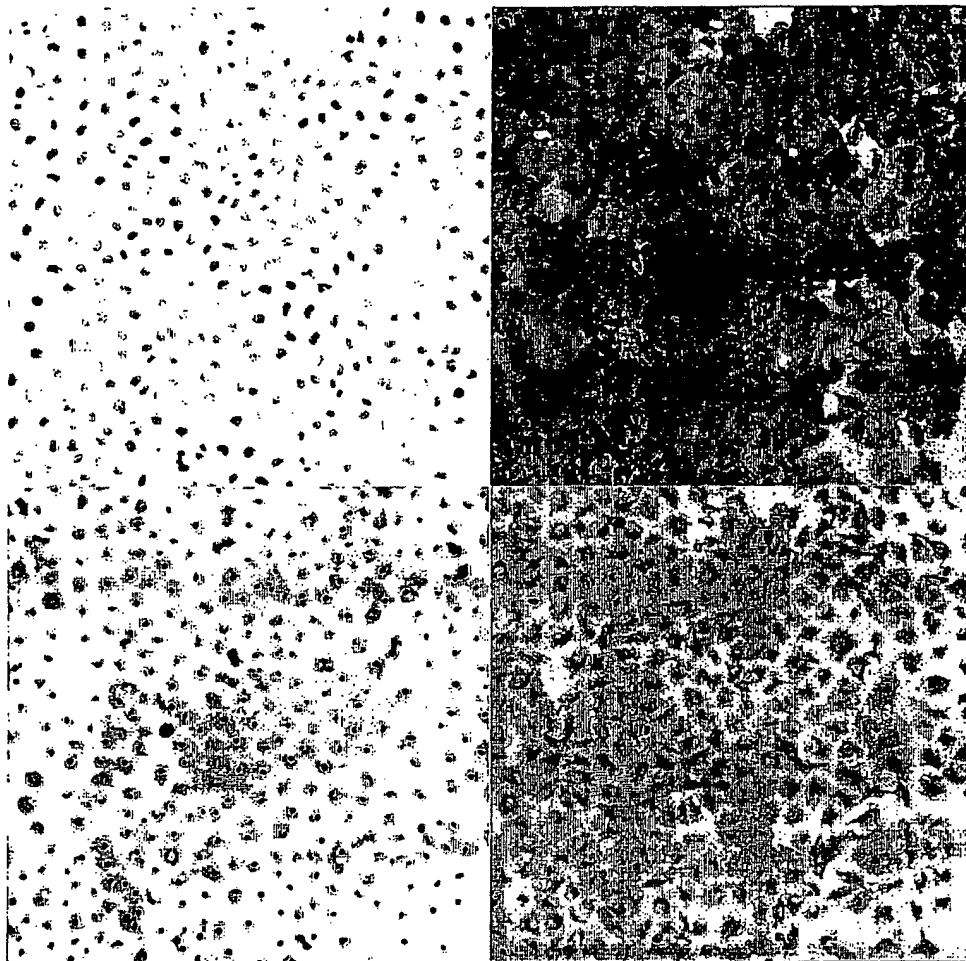


FIG. 13

28/37

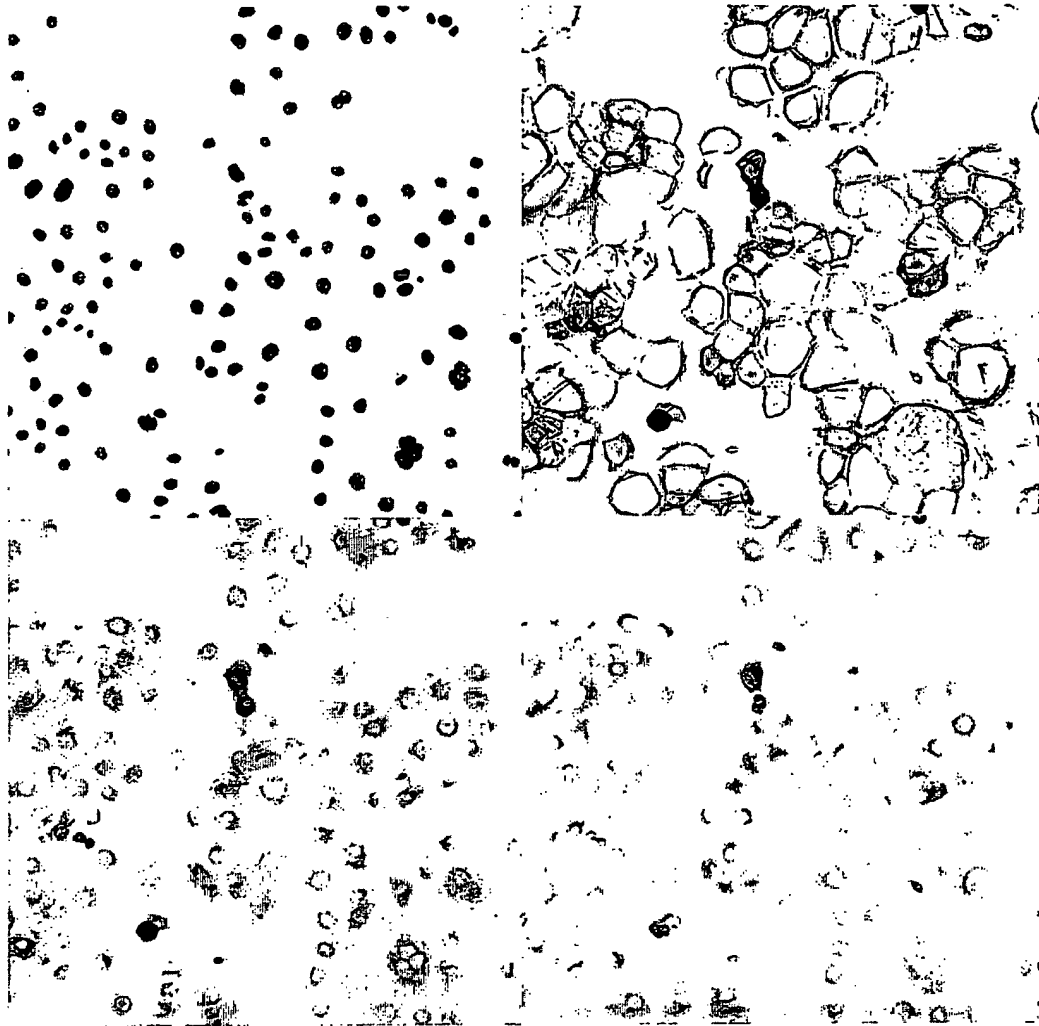


FIG. 14

29/37

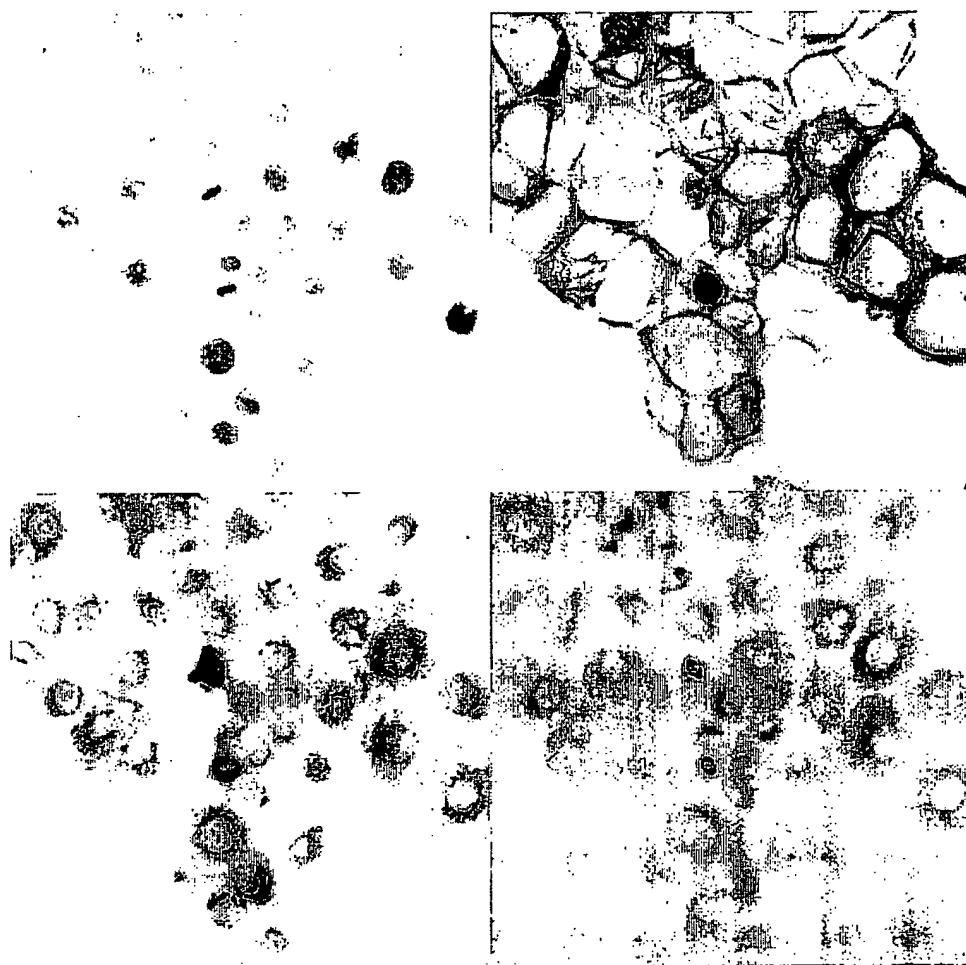


FIG. 15

30/37

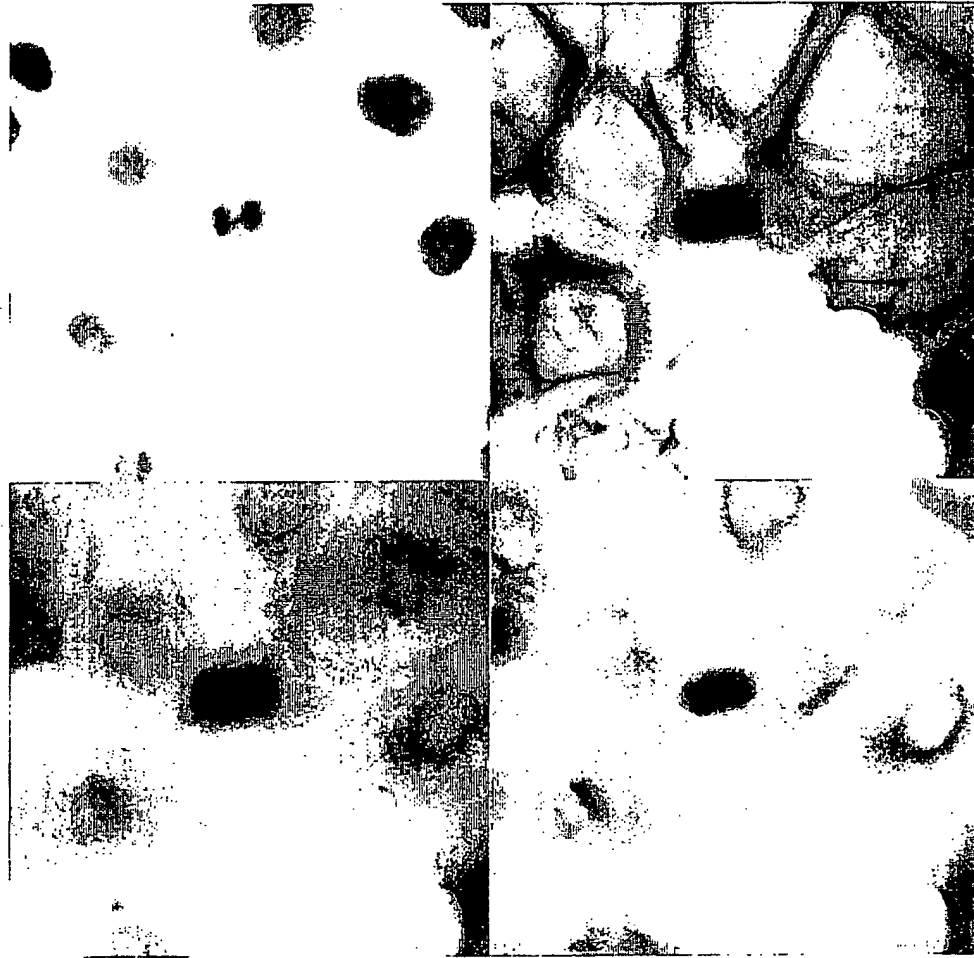


FIG. 16

31/37

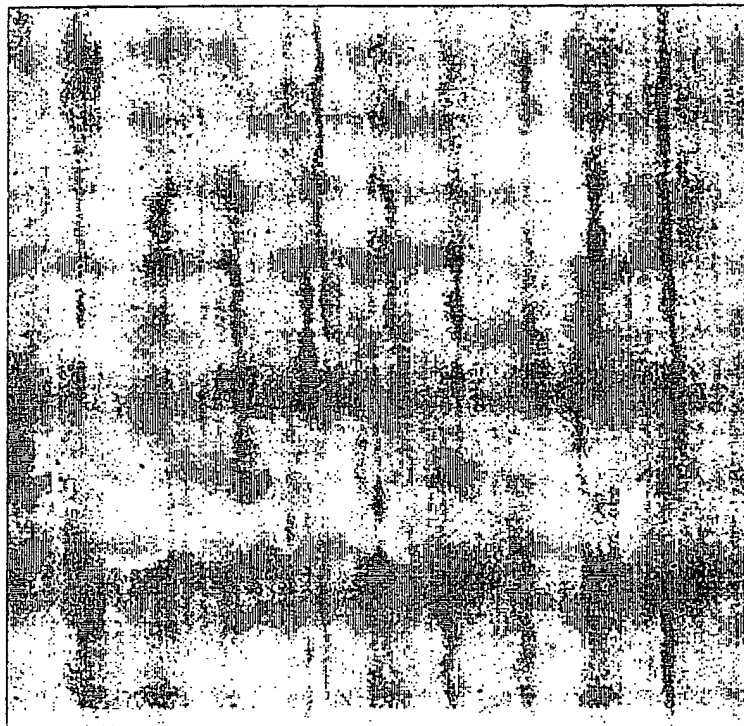


FIG. 17

32/37

Conversion of morphometric parameters into nucleic acid code  
and clustering of the resulting sequences using Neighbor  
Joining method.

Compound:	Measurements																			
	Count	Area	Perimeter	Length	Breadth	Fiber length	Fiber breadth	Shape factor	Ell. form factor	Inner radius	Outer radius	Mean radius	Equiv. radius	Equiv. sphere vol.	Equiv. prolate vol.	Equiv. oblate vol.	Equiv. sphere surface area	Average gray value	Total gray value	Optical density
Control	t	t	t	t	t	t	t	t	t	t	t	t	t	t	t	t	t	t	t	t
Taxol	a	t	t	t	t	t	t	t	a	t	t	t	t	t	t	t	t	t	t	t
CD	c	a	a	a	t	a	t	t	c	a	a	a	a	a	a	a	a	t	a	a
Nocodazol	c	t	t	t	t	t	t	t	t	t	t	t	t	t	t	t	t	t	t	t
Staurosporine	g	g	c	a	a	t	a	a	t	g	a	a	a	t	g	g	g	a	a	t
Vinblastine	c	t	t	t	t	t	t	t	t	t	t	t	t	t	t	t	t	g	t	t
Hydroxyurea	g	t	t	t	t	t	t	g	t	t	t	t	t	t	t	t	t	t	c	t

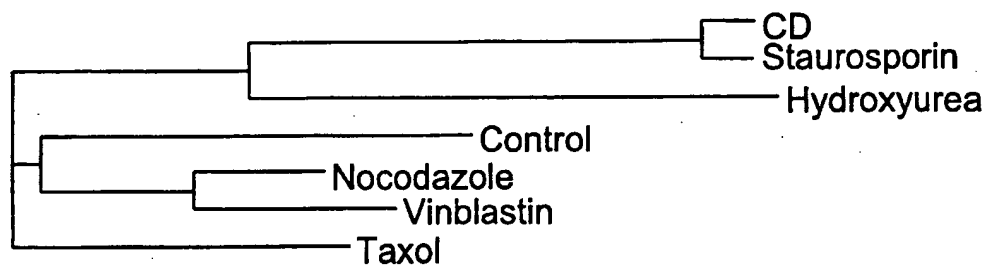


FIG. 18

33/37

Conversion of morphometric parameters into amino acid codes  
and clustering of the resulting sequences using Neighbor  
Joining method.

	Count	Area	Perimeter	Length	Breadth	Fiber length	Fiber breadth	Shape factor	Ell. form factor	Inner radius	Outer radius	Mean radius	Equiv. radius	Equiv. sphere vol.	Equiv. prolate vol.	Equiv. oblate vol.	Equiv. sphere surface area	Average gray value	Total gray value	Optical density	Radial dispersion	Texture Difference Moment	EFA Harmonic 2, Semi-Major Axis	EFA Harmonic 2, Semi-Minor Axis	EFA Harmonic 2, Semi-Major A
Control	H	P	T	T	N	S	D	W	E	S	T	T	T	F	C	C	P	P	M	C	T	G	T	T	Y
Taxol	G	F	M	M	P	M	P	H	G	S	M	M	W	C	F	P	F	R	C	M	M	H	M	P	S
CD	F	G	G	G	M	G	M	K	A	G	G	G	G	G	G	G	G	H	G	G	G	M	G	V	H
Nocodozol	W	F	M	M	W	M	P	T	R	S	M	M	M	F	M	W	F	M	M	R	M	M	M	F	G
Staurosporine	N	V	A	G	G	M	G	G	Y	V	G	G	G	M	V	V	V	G	G	H	G	M	G	G	V
Vinblastine	F	W	W	M	W	W	C	W	D	S	M	W	W	M	M	M	W	M	V	E	M	M	M	F	P
Hydroxyurea	S	H	H	H	H	H	H	V	H	H	H	H	H	H	H	H	H	H	H	A	H	G	H	H	D

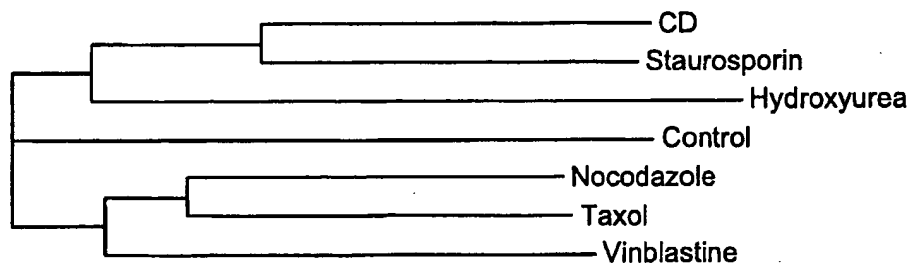


FIG. 19

34/37

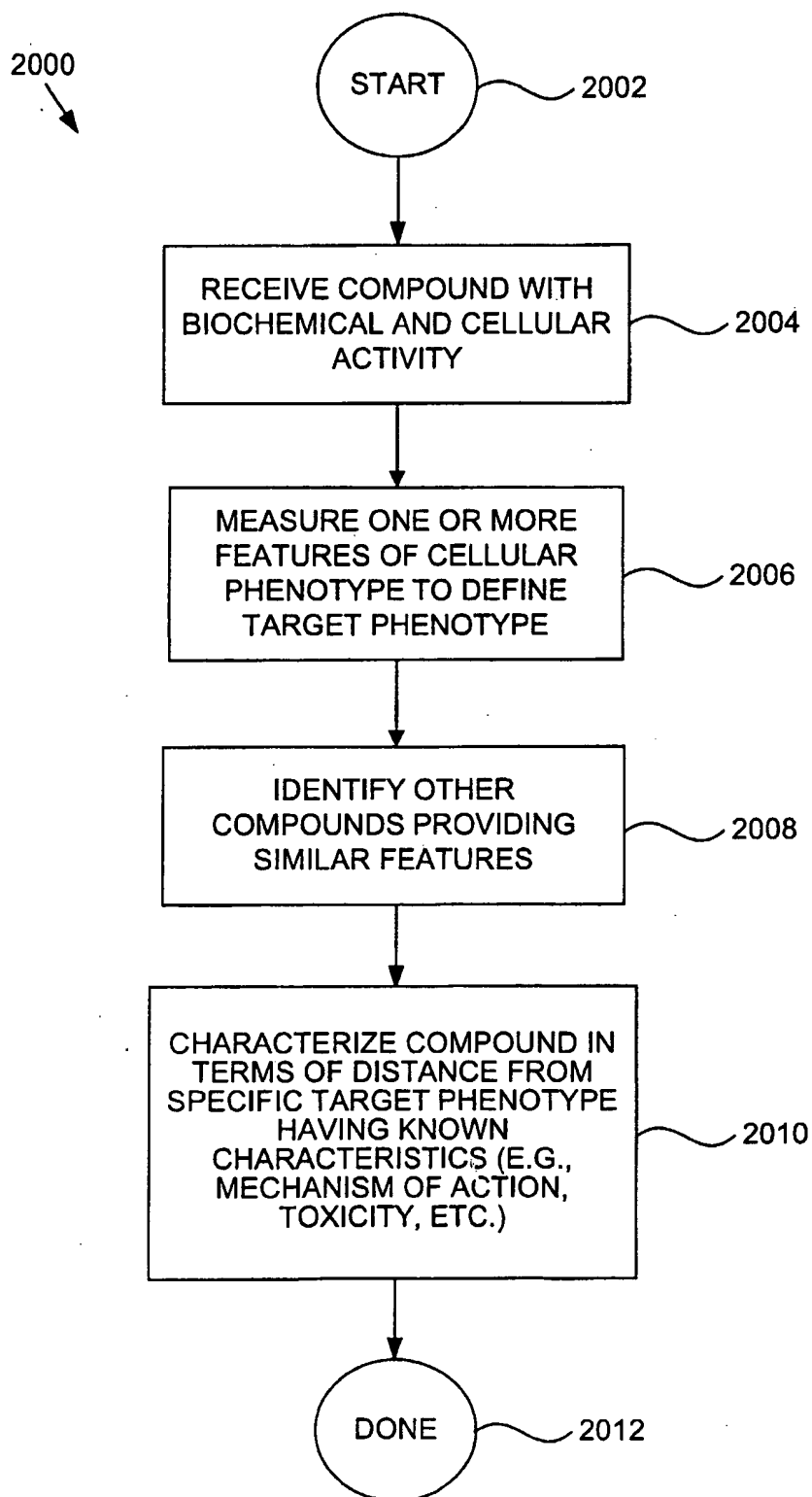


FIG. 20



35/37

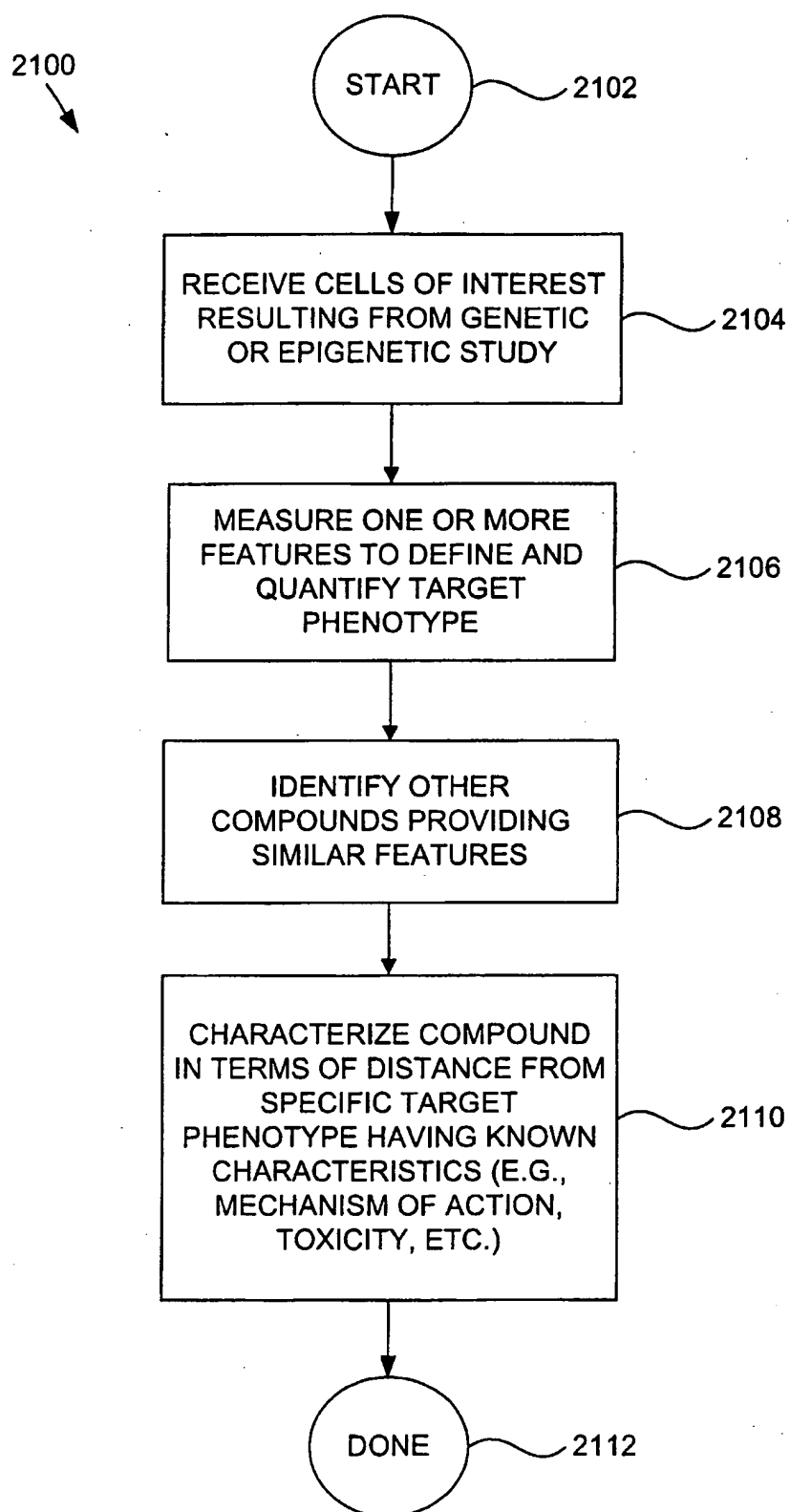
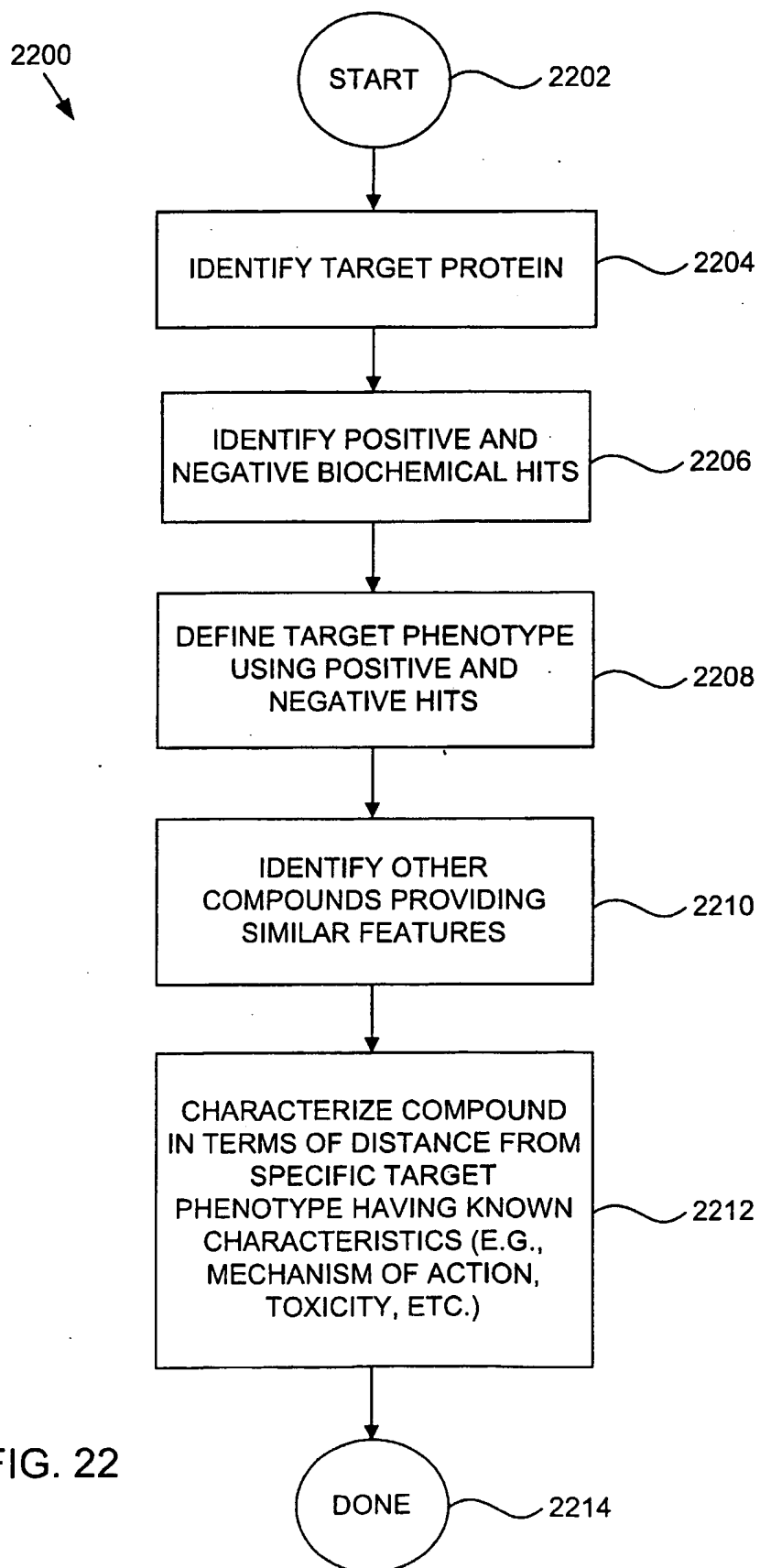


FIG. 21

36/37



37/37

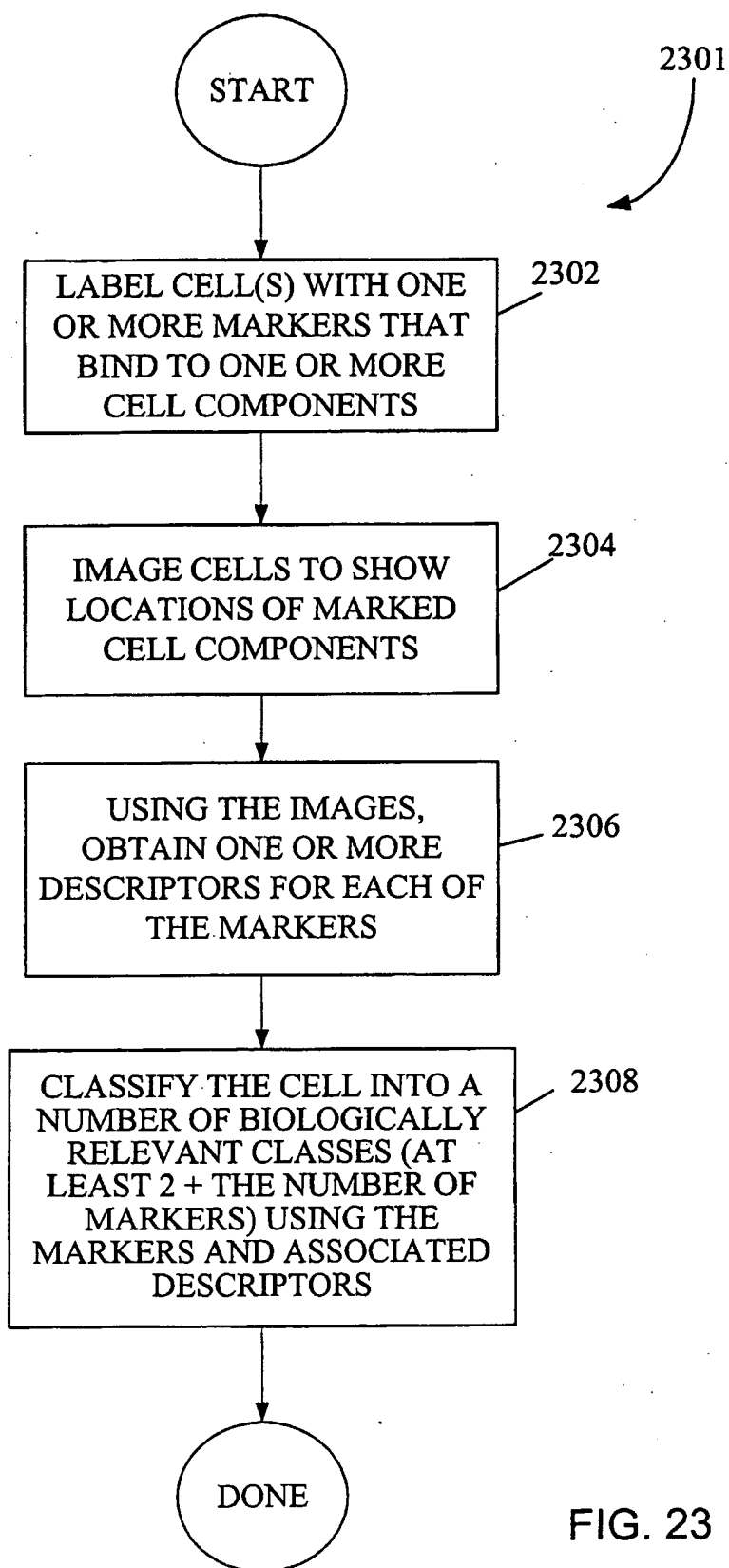


FIG. 23

(19) World Intellectual Property Organization  
International Bureau



(43) International Publication Date  
1 November 2001 (01.11.2001)

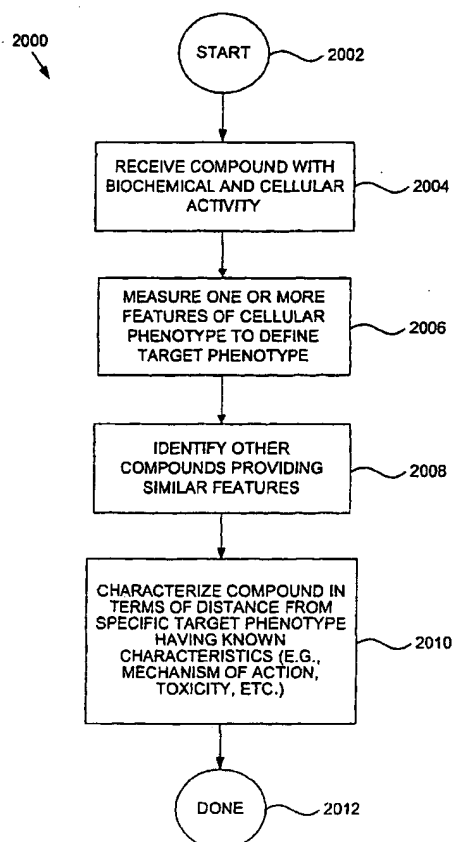
PCT

(10) International Publication Number  
**WO 01/081895 A3**

- (51) International Patent Classification<sup>7</sup>: **G01N 15/14, G06K 9/00**
- (21) International Application Number: **PCT/US01/13248**
- (22) International Filing Date: **24 April 2001 (24.04.2001)**
- (25) Filing Language: **English**
- (26) Publication Language: **English**
- (30) Priority Data:  
60/199,778 26 April 2000 (26.04.2000) US  
09/790,214 20 February 2001 (20.02.2001) US
- (71) Applicant (for all designated States except US): **CYTOKINETICS, INC.** [US/US]; 280 East Grand Avenue, Suite 2, South San Francisco, CA 94080 (US).
- (72) Inventors; and  
(75) Inventors/Applicants (for US only): **OESTREICHER, Donald, R.** [US/US]; 904 Old Town Court, Cupertino, CA 95014-4024 (US). **SABRY, James, H.** [CA/US]; 4305 20th Street, San Francisco, CA 94114 (US). **ADAMS, Cynthia, L.** [US/US]; 2409 Cedar Street, Berkeley, CA 94708 (US). **VAISBERG, Eugeni, A.** [US/US]; 647 Pegasus Lane, Foster City, CA 94404 (US). **CROMPTON, Anne, M.** [US/US]; 2 Bellaire Place, San Francisco, CA 94133 (US).
- (74) Agent: **WEAVER, Jeffrey, K.**; Beyer Weaver & Thomas, LLP, P.O. Box 778, Berkeley, CA 94704-0778 (US).
- (81) Designated States (national): AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, BZ, CA, CH, CN, CR, CU, CZ, DE, DK, DM, DZ, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ,

[Continued on next page]

(54) Title: **METHOD AND APPARATUS FOR PREDICTIVE CELLULAR BIOINFORMATICS**



(57) Abstract: Techniques for using information technology in therapeutics or drug discovery. In an exemplary embodiment, techniques for determining information about the properties of substances based upon information about structure of living or non-living cells exposed to substances are provided. A method according to the present invention enables researchers and/or scientists to identify promising candidates in the search for new and better medicines or treatments using, for example, a multiple biological descriptors derived from a single cell component or marker. The method employs image analysis to extract a plurality of features (e.g., cell size, distance between cells, cell population, cell type) from an image acquisition device into the database.

WO 01/081895 A3



NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM,  
TR, TT, TZ, UA, UG, US, UZ, VN, YU, ZA, ZW.

**Published:**

— with international search report

(84) **Designated States (regional):** ARIPO patent (GH, GM, KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE, TR), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GW, ML, MR, NE, SN, TD, TG).

**(88) Date of publication of the international search report:**

13 March 2003

*For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.*

## INTERNATIONAL SEARCH REPORT

Internat Application No  
PCT/US 01/13248

A. CLASSIFICATION OF SUBJECT MATTER  
IPC 7 G01N15/14 G06K9/00

According to International Patent Classification (IPC) or to both national classification and IPC

## B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)  
IPC 7 G01N

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practical, search terms used)

EPO-Internal

## C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
P,X	WO 00 70528 A (CYTOKINETICS INC) 23 November 2000 (2000-11-23) abstract	1,17
X	WO 97 43732 A (ONCOMETRICS IMAGING CORP) 20 November 1997 (1997-11-20)  page 2, line 20 -page 4, line 12  -/--	1,10,12, 13,17, 26,28,29

☒ Further documents are listed in the continuation of box C.

☒ Patent family members are listed in annex.

## \* Special categories of cited documents :

- "A" document defining the general state of the art which is not considered to be of particular relevance
- "E" earlier document but published on or after the international filing date
- "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)
- "O" document referring to an oral disclosure, use, exhibition or other means
- "P" document published prior to the international filing date but later than the priority date claimed

"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention

"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone

"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art.

"&" document member of the same patent family

Date of the actual completion of the international search

21 November 2002

Date of mailing of the international search report

06. 12. 2002

Name and mailing address of the ISA

European Patent Office, P.B. 5818 Patentlaan 2  
NL - 2280 HV Rijswijk  
Tel. (+31-70) 340-2040, Tx. 31 651 epo nl,  
Fax: (+31-70) 340-3016

Authorized officer

Zinngrebe, U

## INTERNATIONAL SEARCH REPORT

Internat Application No

PCT/US 01/13248

## C.(Continuation) DOCUMENTS CONSIDERED TO BE RELEVANT

Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	WO 96 09605 A (NEOPATH INC) 28 March 1996 (1996-03-28)	1,10,12, 13,17, 26,28,29
Y	page 6, line 8-27  page 7, line 21-31 page 8, line 20 -page 10, line 10 page 10, line 34 -page 11, line 1 ----	2-9,11, 14-16, 18-25, 27,30-32
Y	WO 98 38490 A (BIODX INC ;DUNLAY R TERRY (US); GOUGH ALBERT H (US); GIULIANO KENN) 3 September 1998 (1998-09-03) cited in the application page 10, line 5-20 page 11, line 5-19 page 32, line 1-19 page 44; example 1 page 76, line 15 -page 86 ----	2-9,11, 14-16, 18-25, 27,30-32
X	WO 98 45704 A (TULLIN SOEREN ;KASPER ALMHOLT (DK); NOVONORDISK AS (DK); SCUDDER K) 15 October 1998 (1998-10-15) abstract ----	33,48
X	US 5 326 691 A (HOZIER JOHN) 5 July 1994 (1994-07-05) column 3, line 10 -column 4, line 4 column 17, line 13 -column 19, line 15 ----	33
X	WO 93 21511 A (COMBACT IMAGING SYSTEMS LTD ;UNIV RAMOT (IL); SCHORR KON BEN (GB);) 28 October 1993 (1993-10-28) abstract ----	33
A	WO 99 39184 A (HARTMANN THOMAS ;RIBOZYME PHARM INC (US)) 5 August 1999 (1999-08-05) abstract ----	33,48
A	WO 99 05323 A (AFFYMETRIX INC) 4 February 1999 (1999-02-04) abstract -----	33,48

# INTERNATIONAL SEARCH REPORT

Int'l application No.  
PCT/US 01/13248

## Box I Observations where certain claims were found unsearchable (Continuation of Item 1 of first sheet)

This International Search Report has not been established in respect of certain claims under Article 17(2)(a) for the following reasons:

1. ☐ Claims Nos.:  
because they relate to subject matter not required to be searched by this Authority, namely:
2. ☐ Claims Nos.:  
because they relate to parts of the International Application that do not comply with the prescribed requirements to such an extent that no meaningful International Search can be carried out, specifically:
3. ☐ Claims Nos.:  
because they are dependent claims and are not drafted in accordance with the second and third sentences of Rule 6.4(a).

## Box II Observations where unity of invention is lacking (Continuation of Item 2 of first sheet)

This International Searching Authority found multiple inventions in this international application, as follows:

see additional sheet

1. ☒ As all required additional search fees were timely paid by the applicant, this International Search Report covers all searchable claims.
2. ☐ As all searchable claims could be searched without effort justifying an additional fee, this Authority did not invite payment of any additional fee.
3. ☐ As only some of the required additional search fees were timely paid by the applicant, this International Search Report covers only those claims for which fees were paid, specifically claims Nos.:
4. ☐ No required additional search fees were timely paid by the applicant. Consequently, this International Search Report is restricted to the invention first mentioned in the claims; it is covered by claims Nos.:

Remark on Protest

- ☐ The additional search fees were accompanied by the applicant's protest.
- ☒ No protest accompanied the payment of additional search fees.



FURTHER INFORMATION CONTINUED FROM PCT/ISA/ 210

This International Searching Authority found multiple (groups of) inventions in this international application, as follows:

1. Claims: 1-32

Classifying cells in three or more biological classes using morphological or statistical descriptors

2. Claims: 33-62

Characterising cellular activity of a compound by phenomenological analysis and comparison with reference cells

# INTERNATIONAL SEARCH REPORT

Information on patent family members

International Application No

PCT/US 01/13248

Patent document cited in search report		Publication date	Patent family member(s)	Publication date
WO 0070528	A	23-11-2000	AU 4847100 A EP 1188139 A2 WO 0070528 A2	05-12-2000 20-03-2002 23-11-2000
WO 9743732	A	20-11-1997	US 5889881 A AU 2377997 A WO 9743732 A1 EP 0901665 A1 JP 2000510266 T US 6026174 A	30-03-1999 05-12-1997 20-11-1997 17-03-1999 08-08-2000 15-02-2000
WO 9609605	A	28-03-1996	US 5978497 A AU 3629795 A WO 9609605 A1 US 6137899 A US 6134354 A	02-11-1999 09-04-1996 28-03-1996 24-10-2000 17-10-2000
WO 9838490	A	03-09-1998	US 5989835 A US 6103479 A AU 730100 B2 AU 6667898 A EP 0983498 A1 JP 2000509827 T WO 9838490 A1 US 6416959 B1 AU 734704 B2 AU 3297197 A EP 0912892 A1 JP 2000512009 T	23-11-1999 15-08-2000 22-02-2001 18-09-1998 08-03-2000 02-08-2000 03-09-1998 09-07-2002 21-06-2001 05-01-1998 06-05-1999 12-09-2000
WO 9845704	A	15-10-1998	AT 215227 T AU 6820998 A DE 69804446 D1 DE 69804446 T2 WO 9845704 A2 DK 986753 T3 EP 1199564 A2 EP 0986753 A2 ES 2173573 T3 JP 2001522454 T	15-04-2002 30-10-1998 02-05-2002 07-11-2002 15-10-1998 22-07-2002 24-04-2002 22-03-2000 16-10-2002 13-11-2001
US 5326691	A	05-07-1994	AU 3134893 A WO 9310259 A1 US 5563060 A	15-06-1993 27-05-1993 08-10-1996
WO 9321511	A	28-10-1993	AT 182211 T AU 3899393 A CA 2117777 A1 DE 69325652 D1 DE 69325652 T2 EP 0635126 A1 WO 9321511 A1 JP 8511676 T RU 2126962 C1	15-07-1999 18-11-1993 28-10-1993 19-08-1999 28-10-1999 25-01-1995 28-10-1993 10-12-1996 27-02-1999
WO 9939184	A	05-08-1999	AU 2241899 A WO 9939184 A1	16-08-1999 05-08-1999

# INTERNATIONAL SEARCH REPORT

Information on patent family members

International Application No

PCT/US 01/13248

Patent document cited in search report		Publication date	Patent family member(s)	Publication date
WO 9905323	A	04-02-1999	EP 1002264 A2	24-05-2000
			EP 1009861 A1	21-06-2000
			EP 0998697 A1	10-05-2000
			EP 1007737 A1	14-06-2000
			JP 2001511546 T	14-08-2001
			JP 2001515234 T	18-09-2001
			JP 2001511550 T	14-08-2001
			JP 2001511529 T	14-08-2001
			WO 9905323 A1	04-02-1999
			WO 9905574 A1	04-02-1999
			WO 9905324 A1	04-02-1999
			WO 9905591 A2	04-02-1999
			US 6229911 B1	08-05-2001
			US 6188783 B1	13-02-2001
			US 6308170 B1	23-10-2001
			US 2001018642 A1	30-08-2001
			US 2002012456 A1	31-01-2002
			US 2002062319 A1	23-05-2002
			US 2002150932 A1	17-10-2002

**This Page is Inserted by IFW Indexing and Scanning  
Operations and is not part of the Official Record**

**BEST AVAILABLE IMAGES**

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images include but are not limited to the items checked:

- ☐ **BLACK BORDERS**
- ☐ **IMAGE CUT OFF AT TOP, BOTTOM OR SIDES**
- ☐ **FADED TEXT OR DRAWING**
- ☐ **BLURRED OR ILLEGIBLE TEXT OR DRAWING**
- ☐ **SKWEDED/SLANTED IMAGES**
- ☒ **COLOR OR BLACK AND WHITE PHOTOGRAPHS**
- ☐ **GRAY SCALE DOCUMENTS**
- ☐ **LINES OR MARKS ON ORIGINAL DOCUMENT**
- ☐ **REFERENCE(S) OR EXHIBIT(S) SUBMITTED ARE POOR QUALITY**
- ☐ **OTHER:** \_\_\_\_\_

**IMAGES ARE BEST AVAILABLE COPY.**

**As rescanning these documents will not correct the image problems checked, please do not report these problems to the IFW Image Problem Mailbox.**